

Project II

Group 01 - M. Carli , C. Meyer, O. Tausendschön & B. Takács

2023-11-15

Contents

Report	2
Introduction / Abstract	2
Our Aims	2
Preparation	2
Scraping the Data	2
General Analysis	3
Text Analysis	5
Sentiment analysis	7
Emotion Analysis	8
Topic Analysis	10
Modeling approaches	11
Summary	13
Appendix	14
Scraping the data	14
General analysis	18
Text Analysis	22
Sentiment analysis	28
Emotion Analysis	31
Topic Analysis	36
Modeling approaches	40

Report

Introduction / Abstract

In this Report, we take a deep dive into the world of customer reviews to uncover trends and insights related to a selected product on Amazon.com. With a foundation of about 500 text-based evaluations, the product that we have carefully picked from a product category of our choosing offers big variety of thoughts and comments.

The R code of the entire process is available as an appendix, revealing the details of our analysis procedure. Additionally, our presentation summarises the most significant discoveries, providing a quick overview of the brand's online presence and consumer attitude. The gathered dataset, in .rda format provides the basis for our investigation and creates a clear picture of the brand's interaction with its clientele.

Our Aims

Our goal is to analyse and extract useful data from the collection of customer reviews. In doing this, we want to provide basic answers that explain the brand's effectiveness and consumer impression.

In order to get a good impression of customer's opinions, we will analyse common and frequent terms that appear in the reviews of this data. We will also explore temporal dimension, tracking indicators as they change over time to provide insights how the product or the opinions evolve over time. We will be examining the tone of the text to determine the emotional undertones that are typical of customer feedback. Finally, we will interpret the subjects that predominate in customer conversations, revealing the problems and elements that are most important to the clientele.

Our ultimate objective is to offer the brand with insights that can be put into practise as we go through these elements. If the brand's goals are a high star rating and happy consumers, our study will help them in archieving that.

Preparation

Before analysing the data, we actually have to get the data. We did this by using a basic webscraper to extract important information on the prodct's review page.

Scraping the Data

To scrape reviews which we can then analyze we start by installing necessary and helpful packages. We could then start with tasks such as HTML parsing and mimicking a browser with a machine.

However, we quickly ran into problems as Amazon makes an effort to prevent scraping, particularly when many pages are being scraped. We started by trying the code with different products and pages, experimenting with different HTML nodes and xpaths. All this was done with the R-package **Rvest**. In the end, we created custom headers, a random time out to to mimic human behaviour as well as to avoid being "too aggressive" and imitated request headers. We managed to scrape 100 reviews without filtering for specific terms. More reviews where not attainable, since Amazon, as part of their web-scraping prevention, limits the number of pages of reviews to be viewed to 10, and the number of reviews per page also to 10. To increase the dataset size, we filtered for all five possible star ratings and obtained the 100 available reviews per star rating, resulting in a total of 500 (5x100) reviews. However, this comes with the big downside of introducing a bias. We nevertheless

decided to continue as one purpose of this project is to carry out a complete analysis and larger, although biased, dataset allows for interesting insights into the different wordings, sentiments and emotions associated with different star ratings.

We decided to extract information on the title of the reviews, the review content, the review date, whether the review is verified or not, how many people found it helpful, and the star rating. Our dataset thus contains 500 observations with 6 variables.

The variables are the following:

- **review_title**: This is the title of the review
- **review_text**: This is the review itself.
- **formatted_date**: This is the date when the review was published.
- **verified**: This indicated whether a review is made by an officially verified buyer or not.
- **N_helpful**: This is the count of people that marked the review as being helpful.
- **star_rating_num**: This is the score of the review. The maximum is five and minimum is zero. A higher number corresponds to a more positive opinion about the product.

The technical details of the data analysis are provided in the appendix, accompanied by the corresponding R code.

General Analysis

To get an idea of the data we are working with, We want to start by taking a look at some general features of our review data. We will keep this section short as it only serves to get a basic understanding of the data.

Let us start by taking a look at **when the reviews were made**: The time frame of our reviews spans from 06th of June, 2020, when the first review was published until the 17th of November, 2023, where the most recent one was written. This shows and ensures that we consider reviews in a long time span, considering potential product updates, relaunches and alike.

It is also interesting to get more of an insight into the **general rating behavior**. The **worst rating** is 1 stars, the **highest one** is 5 stars and the **average rating** amounts to 3 stars. This can be shown by using a histogram of the distribution of star ratings. Note that the equal distribution is due to the way our webscraper works.

Another visualization of the **distribution of star ratings** might be helpful to see how the product has performed in general. Thus, we present a histogram:

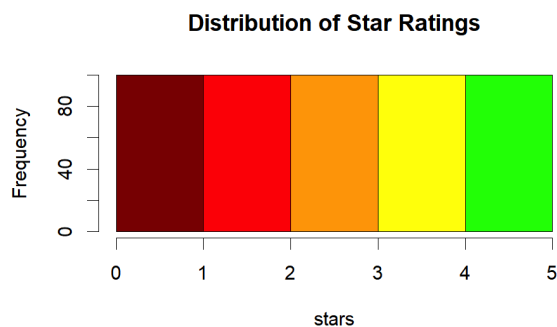


Figure 1: Distribution of star ratings.

We can also visualize the average star rating across time. This gives us a good guess of how the quality of the product might deteriorate or improve.

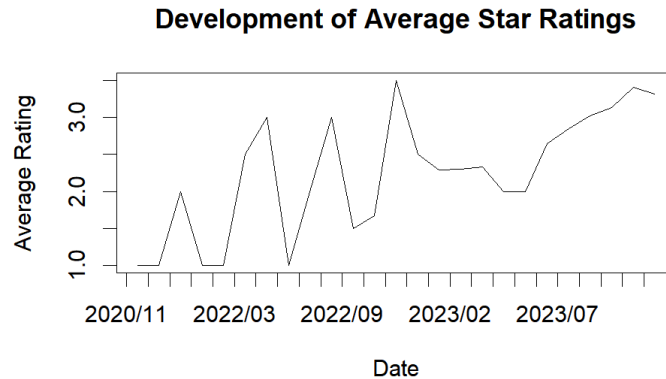


Figure 2: Development of average star ratings over time.

It seems as the average rating improves in the long term, which is a good sign. Let's see if there is any correlation with **the number of reviews published across our time period**.

Also, the average rating seems to improve already after the first reviews and stays relatively constant throughout the time period where reviews were published.

It might also be interesting to take a look at **the number of reviews published across our time period** and see if and how this corresponds to the development of the star ratings.

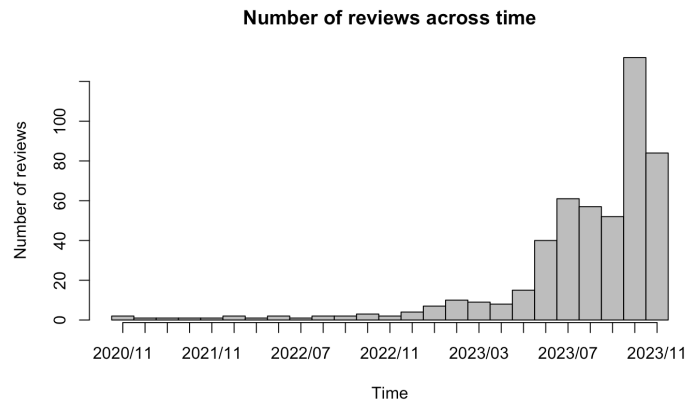


Figure 3: Number of reviews across time.

The plot suggest that most of the reviews were published recently.

When comparing the two plots we can see that overall the review-situation improved over time: Both more reviews and better reviews were published in the long run.

Furthermore, we might want to take a look at **helpfulness of the reviews**. We can see that the reviews were rated as helpful by between 0 and 50 people, with an average of 3.47 and a median of 0 as well. This indicates skewedness and strong outliers, so we want to take a closer look:

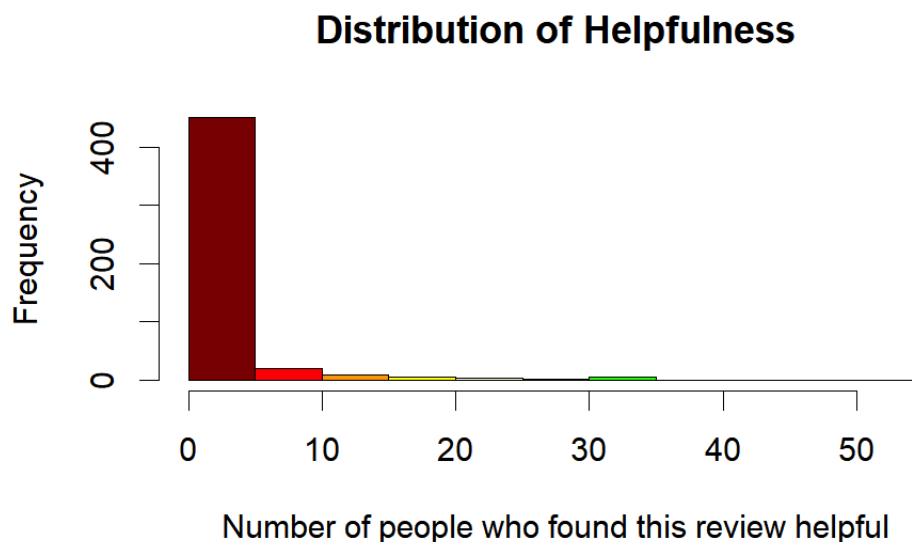


Figure 4: Frequency of helpful reviews.

This supports our hypothesis that the vast majority of reviews are rated as helpful by very few to none people. Only very few people attained over 15 “helpfulness”-votes. We can see the “Winner takes it all” principle at play. Overall a third of the reviews were perceived as helpful, meaning that at least 1 customer rated them as such.

Text Analysis

This section focuses on text analysis, aiming to uncover patterns, relationships, and insights embedded in the textual feedback provided by customers.

To kick off our text analysis, we start by taking a look at review length. We can see that the average length of a review in our dataset is about 70 words. But what might be more interesting is the correlation between the length of a reviews and its rating. The correlation is about -0.3. What this reveals is that long reviews tend to be more negative or that customers who have a negative opinion about a product tend to share more information to explain or complain.

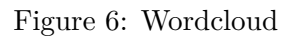
The correlation between the length of a review and its helpfulness is at about 0.63. This leads to the conclusion that longer reviews tend to be rated as more helpful as they for example share more information and details that can be useful for other potential customers.

Now let us dive even deeper into text analysis than merely the word count. Let’s have a look at **what people actually write in their reviews**. To do this, we cleaned the reviews by removing special characters, nubers, punctuation, extra white spaces, common stopwords and reduced all words to their stem. After unsurprisingly seeing words like Phone, iphone, apple as the most frequent, we removed these and constructed a wordcloud (see Figure 5).

This reveals a much more interesting picture. We chose to show all terms that are mentioned at least 5 times in the 100 reviews in the word cloud. The by far biggest and thus most frequently



Let us now do the same for the titles of the reviews and see if there are differences:



6

Sentiment analysis

Before we dive into applying sentiment analysis on our dataset of reviews, we tried to recall the basics of the Sentiment package and its meaning in R. Essentially sentiment analysis works by differentiating between words depending on the sentiment that is attached to them. This is conducted via dictionaries consisting of lists of positive vs. negative words, or lists of more diverse emotions. Packages such as `sentimentr` in R work by scanning the text to see if words in the text match with any dictionary entries. The words are then assigned a value (>0 if the word is located on the positive list, <0 if it is on the negative one) - all values are added together and the average sentiment is determined. Important note: The packages take the words before and after a term into account in order to assess its classification. This way valence shifters or negations can be included. Sentiment analysis can be applied in a number of fields and situations. Its use-cases range from social media monitoring, political campaigns, PR to market research.

We now want to start our first computations in the field of sentiment analysis to get a better picture about the general tonality and sentiment of the reviews we are examining. In this we will differentiate between title and text, look at the correlation between them and their general behavior.

The computations reveal that in the review texts, sentiment ranges from -0.54 to 1.42 and averages at around 0.25. For the titles, the lowest sentiment is -1.44, the highest is 1.23 and the average is around 0.35. Both average sentiments are positive, which is good news for the product. Anything else would be highly unusual since the sentiment should correspond to star ratings and since there are no significant bad ratings, this would need further investigation. We also prove that the sentiment of title and review text are correlated since the p-value is extremely low.

Now, similar to the way we wanted to check if the rating behavior changed over time we want to take a look at the **development of sentiment over time**:

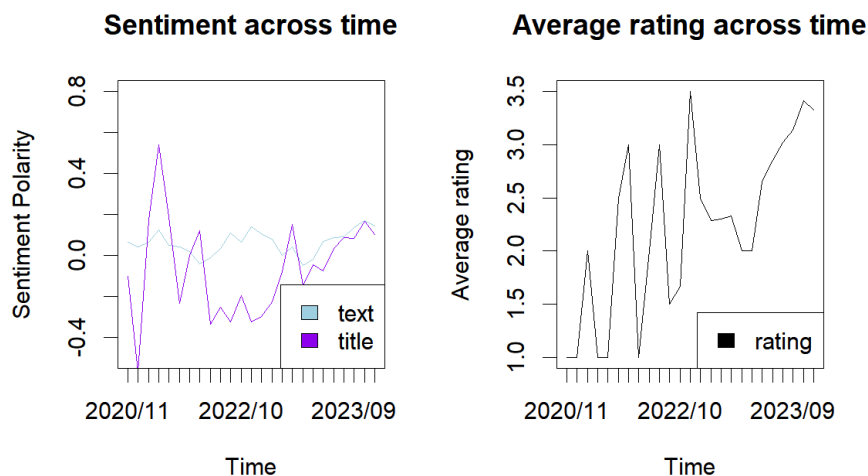


Figure 7: Sentiments

When comparing the plots of how sentiment (in both title and text) developed, we can see a clear correlation with average rating. This is good to see as it confirms that users who put a better rating put more positive emotions into the review.

Let us now analyse how sentiment and other variables interact in more detail:

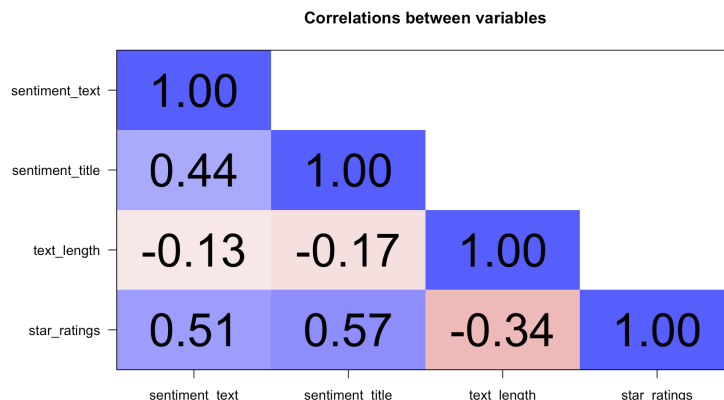


Figure 8: Corplot

In this correlation plot or matrix we can see a rather strong positive correlation between sentiment in title and sentiment in text - this we have already found out. Additionally we can see slight negative correlations between text length and sentiment. This further supports our findings that longer reviews will have a worse rating. Star ratings thus negatively correlate with text length, but are positively impacted by the sentiment score. All in all, we are glad to see this as it confirms our previous assumptions.

Emotion Analysis

With emotion analysis we aim to get more detailed insights into the emotions that are expressed in reviews by differentiating not only between positively and negatively connotated terms, but a more developed spectrum of emotions. As a basis, Plutchik's wheel of emotion is used.

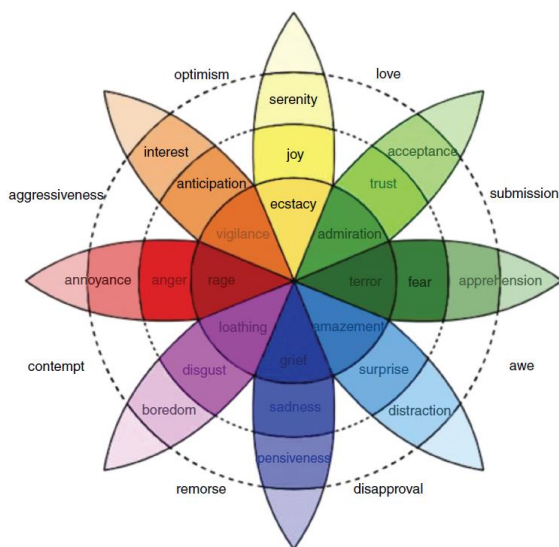


Figure 9: Wheel of Emotion

In order to obtain insights into the emotions in the review texts, we categorize the extracted emotions into the categories of the Wheel of Emotion and examine the degree of average emotion.

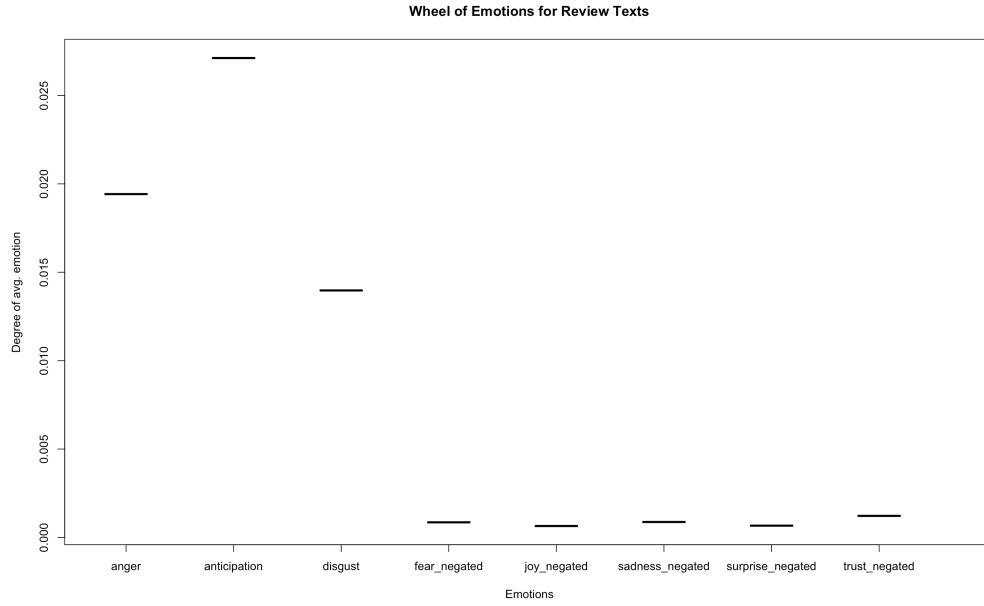


Figure 10: Categorization of review text according to Wheel of Emotion.

Furthermore, in order to determine the relevance of the individual emotions, we examine their prevalence among the review texts:

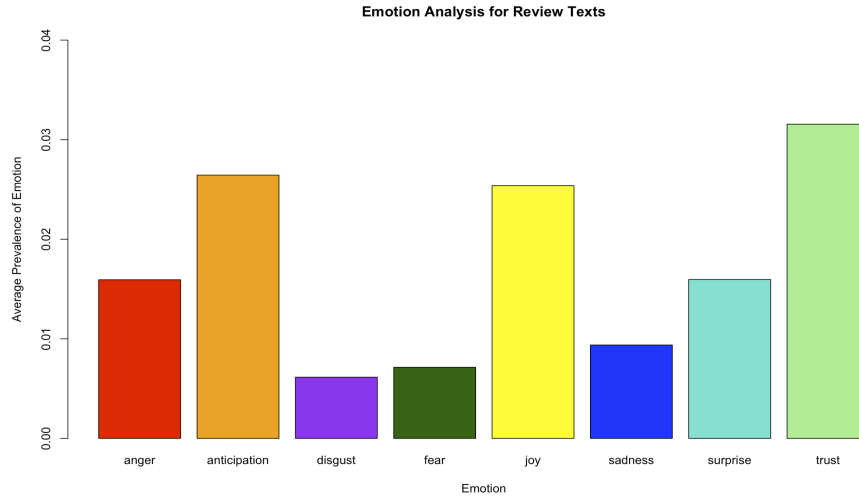


Figure 11: Prevalence of emotions.

The bar chart above shows the determined prevalence of the individual emotions. The higher the bar, the higher the relevance of the corresponding emotion. Based on the plot, we can infer the most dominant emotions in the review texts: the latter are anticipation (orange bar), joy (yellow bar) and trust (light green bar). However, also anger plays an important role in the reviews as well as the emotion of surprise.

To increase our understanding of the emotions expressed in the reviews, we chose to examine the correlation between the individual emotions to uncover possible connections and insights into how the different emotions might go hand in hand with another in the reviews we are examining. The therefore obtained results are shown in a correlation matrix.

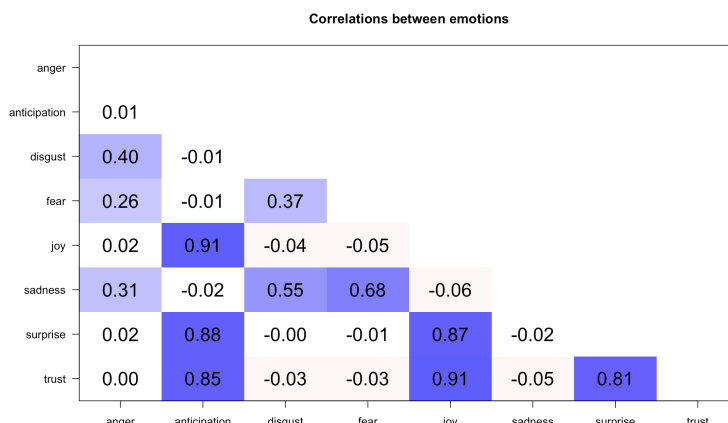


Figure 12: Correlation among different emotions.

The results reveal a high positive correlation between joy and anticipation, joy and trust, anticipation and surprise and trust as well as between surprise and anticipation. No noteworthy negative correlation among emotions seems to be present.

Topic Analysis

After getting a feel for the sentiment and emotions present in the reviews for our production under examination, we next set out to analyse the topic of the individual reviews, trying to identify the latent topics among them and to determine what people are talking about. To do this efficiently, we will make use of topic modelling, which refers to methods that aim to identify topics inside a text corpus.

In our case, we use Latent Dirichlet Allocation (LDA), which is a popular probabilistic model in topic modelling. LDA assumes that documents are a mixture of topics and that each word in a document can be attributed to one of the document's topics. To identify latent topics, the model analyses the distribution of words across a collection of documents.

Process-wise after pre-processing, a document term matrix is created, the number of topics to differentiate between is decided upon, we then check for convergence, estimate the model & can then interpret our findings. The number of topics, which is a key hyperparameter for LDA models, was determined semi-automatically, testing different combinations and choosing the model with the lowest Akaike Information Criterion (AIC). The results suggested that three topics are the optimal choice in our case. (*Note: The detailed procedure can be found in the Appendix.*)

Applying our LDA model, with the number of topics to be determined to three, we identified the following topics:

- **Quality:** People seem to be talking about the quality of the product and the purchasing process. Words like “good”, “work” and “excel” are associated with it.
- **Hardware:** In this case, people seem to be talking about technical aspects, mentioning words “battery”, “scratch” and “speaker”.

- Value: For this topic, people seem to be addressing the value of the purchased product, which is in our case a refurbished/renewed product. Mentioning word like “great”, “new”, “life”, “condit” and “worth”.

Lastly, we examine the relative importance of the individual topics in our review, which is shown in the chart.

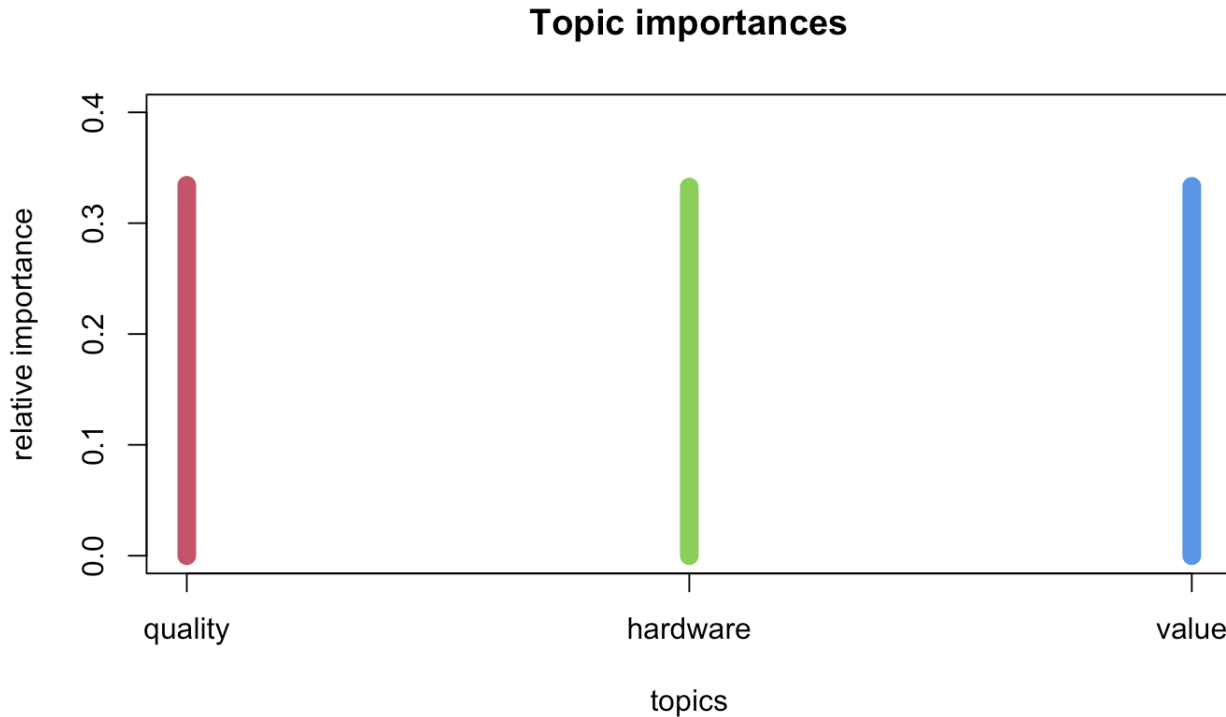


Figure 13: Relevance of individual topics.

Interestingly, all identified topics seem to have the same average relative importance. It is also worth to note that the lowest relative importance that a topic has in our entire dataset was 0.28 (for the topic “quality”) and the maximum value was .39 for “hardware”. The 1st quantile is .32 for all topics, and “quality” has the highest 3rd quantile value at under .35. This means that all three topics are approximately evenly present in the overwhelming reviews, and therefore are not suitable predictors for star rating.

However, one must critically note that this could also be the case due to the dataset generation process and the hereby introduced bias. Further analysis and more data would be necessary to achieve a clear differentiation.

Modeling approaches

After conducting the above analyses, we move on to modeling the star ratings of our product using different features that we have derived in the previous chapters. Our aim with the modeling is to better understand which factors are a significant predictor of ratings.

Star rating As our first target variable is star rating, for which we build linear regression models. In our first model, we used date as the only predictor, and found that there is only a significant difference in ratings in October compared to our baseline of January, with an increase of 1.06 from 2.29.

Next, we used review length as a predictor and found that it negatively impacts ratings with a statistical significance. We kept review length as a predictor in all of our following models.

Afterwards, we found that a higher sentiment score of both the review text and title has a significant, positive effect on ratings.

Next, we looked at emotion: out of the eight emotions, only joy and disgust deemed significant, unsurprisingly with the former having a positive and the latter having a negative effect on ratings. However, with the combination of sentiment scores and emotions, both the title and text sentiments were significant, but none of the emotions.

Adding back dates to the previous model, we find that none of them are significant. We therefore arrived at the conclusion that the best combination of predictors to use are the sentiment scores and review length. To validate this claim, we looked at AIC scores and found that this model had the lowest value, out of the ones created.

However, with stepwise variable selection we could slightly reduce AIC further. This model added the emotions of anticipation, fear and joy to our simplified model. Our model of choice therefore is:

```

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)    3.0743044  0.0740881  41.495  < 2e-16 ***
ave_emotion.anticipation -3.3261650  1.6400355  -2.028   0.0431 *
ave_emotion.fear      -5.7474302  3.3109845  -1.736   0.0832 .
ave_emotion.joy        2.9548995  1.7392872   1.699   0.0900 .
sentiment_title$ave_sentiment 1.3217669  0.1253352  10.546  < 2e-16 ***
sentiment_text$ave_sentiment 1.6223956  0.2164607   7.495 3.10e-13 ***
review_length    -0.0040314  0.0006143  -6.563 1.34e-10 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1.037 on 493 degrees of freedom
Multiple R-squared:  0.4698,    Adjusted R-squared:  0.4633
F-statistic: 72.8 on 6 and 493 DF, p-value: < 2.2e-16

```

Figure 14: Best linear model

Other models Afterwards, we used logistic regression models to predict whether a review will yield 5 stars. We have found that besides the above mentioned predictors, surprise was also significant, negatively impacting the odds of a five star review. No other emotions were significant. To this, by adding the dates as well we reached a model including all our variables. Although no predictors were significant besides the sentiment scores, AIC favoured the full model. Using stepwise variable selection, we excluded some emotions except anger, joy and surprise, further lowering AIC.

For the prediction of whether at least 10 people will find a review helpful, only review length showed to be significant. As opposed to the ratings, sentiment score did not hold predictive power. Using

stepwise variable selection, we landed on a model that had included review length and title sentiment only.

Summary

Overall, this report presents a comprehensive analysis of customer reviews for a selected product on Amazon.com, with the focus of extracting valuable insights to help the brand achieving high star ratings to satisfy customers. The chosen dataset consists of 500 text-based reviews, collected through web scraping and filtering for different star ratings.

Although getting into challenges and not being able to scrape the entire product page, we managed to carefully analyze a big variety of reviews with mixed emotions, topics and ratings. On a general level, we found out the following:

- Reviews span from June 6, 2020, to November 17, 2023, showing a long-time perspective.
- Average star rating is 3, with a histogram showing an equal distribution due to the web scraping method.
- Improvement in average rating over time, with an increase in both the number and quality of reviews.
- Seasonality observed in the number of reviews, with an overall increasing trend

In regard to text analysis, we saw that the term “battery” is used most frequently, indicating that people talk about this often, making it crucial for the company to address various battery related topics. A Sentiment analysis shows positive sentiments in both titles and review texts. These have a correlation to star ratings. Emotion analysis identifies dominant emotions as anticipation, joy, trust, anger, and surprise.

We also employed a Latent Dirichlet Allocation (LDA) to identify three topics: Quality, Hardware, and Value. All topics are evenly present in reviews, but not significant predictors for star ratings. Similarly to the most frequently used terms, this shows what the company should focus as these topics are highly discussed in the reviews.

Finally, we applied linear regression using predictors such as date, review length, sentiment, and emotions to predict star rating. We decided to use sentiment scores, review length as well as anticipation, fear and joy as predictors as these give us the best model fit to predict star ratings.

Appendix

The appendix contains the technical analysis of the available data, including all the corresponding R-code. For the analysis we used the following R-packages:

Scraping the data

To scrape reviews which we can then analyze we start by installing necessary / helpful packages:

```
if (!require("pacman")) install.packages("pacman")
pacman::p_load(rvest,
               polite,
               tidyr,
               dplyr,
               ggplot2,
               tibble,
               purrr)
library(tidyverse)
library(rvest)
library(vctrs)
library(tm)
library(SnowballC)
library(wordcloud)
library(RColorBrewer)
library(ggplot2)
library(udpipe)
library(sentimentr)
library(textcat)
library(pscl)
library(topicmodels)
library(psych)
```

Next up, we start our web-scraping to read information from a chosen website. Here we faced two major difficulties: Amazon blocks attempts of webscraping when iterated over numerous pages to avoid bots - this we countered by making the request with customer headers. Additionally, however - and we did not find a solution to this - only 100 reviews can be obtained via this code since to see further requests Amazon allows only searches for certain terms. Although we could have implemented some terms to search for and add the reviews to our dataset, we did not think that this is a proper solution. Thus we went with only 100 reviews to start with.

```
# Set parameters for scraping reviews
n_reviews <- 100
reviews_per_page <- 10
iters <- ceiling(n_reviews/reviews_per_page)

# Set the base and additional URL parameters
url_p1 <- "https://www.amazon.com/Apple-iPhone-11-128GB-Black/product-reviews/
B07ZPKR714/ref=cm_cr_getr_d_paging_btm_"
url_p2 <- "?ie=UTF8&reviewerType=all_reviews&pageNumber="
```

```

url_p3 <- "&filterByStar="

filters <- c("five_star", "four_star", "three_star", "two_star", "one_star")

# Initialize an empty data frame to store all reviews
reviews_all <- data.frame(NULL)
reviews_scraped <- data.frame(NULL)

# Set a seed for reproducibility
set.seed(1479)

# Loop through the iterations to scrape reviews
for (filter in filters) {
  for (i in 1:iters) {
    # Construct the URL for each iteration
    url <- paste0(url_p1, ifelse(i == 1, "prev_1", paste0("next_", i)), url_p2,
                    i, url_p3, filter)

    # Set custom headers to mimic a browser request
    headers <- c(
      'User-Agent' = 'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36
      (KHTML, like Gecko) Chrome/91.0.4472.124 Safari/537.36',
      'Accept' = 'text/html,application/xhtml+xml,application/xml;q=0.9,
      image/webp,image/apng,*/*;q=0.8'
      # Add more headers as needed
    )

    # Make the request with custom headers
    response <- httr::GET(url, httr::add_headers(.headers=headers))

    # Check if the request was successful (status code 200)
    if (httr::status_code(response) == 200) {
      # Read the HTML content of the page
      page <- read_html(content(response, "text"))

      # Extract Body of 10 reviews per page
      page <- page %>%
        html_elements(
          xpath = "/html/body/div[1]/div[2]/div/div[1]/div/div[1]/div[5]
          /div[3]/div")

      # Extract review information from the page
      ## Review Title
      review_title <- page %>%
        html_elements(xpath = "//*[@data-hook='review-title']/span[2]") %>%
        html_text() %>%
        str_trim() %>%

```

```

    str_remove_all("<.*?>") %>%
    str_replace_all("\\s+", " ")
## Review Text
review_text <- page %>%
  html_elements(xpath = "//*[data-hook='review-body']") %>%
  html_text() %>%
  str_trim() %>%
  str_remove_all("<.*?>") %>%
  str_replace_all("\\s+", " ")
## Review Star Ratings
star_ratings <- page %>%
  html_elements(xpath = "//*[data-hook='review-star-rating']") %>%
  html_text()
## Review Dates
review_dates <- page %>%
  html_elements(xpath = "//*[data-hook='review-date']") %>%
  html_text()

## Review Verified Purchase
reviews <- html_elements(page, xpath = "//*[data-hook='review']")
review_verified <- data.frame(Verified = rep(FALSE, length(reviews)))
j <- 1
for (element in reviews)
{if (str_detect(str_squish(html_text(element)),
                    fixed("Verified Purchase"))){
  review_verified[j, ] <- TRUE
}
  j <- j+1}

## Review n-helpful
review_n_helpful <- data.frame(N_helpful = rep(NA, length(reviews)))
k <- 1
for (element in reviews) {
  review_helpful <- element %>%
    html_text()

  # Define the pattern to match
  pattern <- "((?:\\d+|One) person|\\d+ people) found this helpful"
  match_result <- str_match(review_helpful, pattern)

  if (!is.na(match_result[1, 2])) {
    n_helpful <- match_result[1, 2]

    if (n_helpful == "One person") {
      review_n_helpful[k, ] <- 1
    } else {
      pattern_2 <- "\\b\\d+\\b"

```



```

        extracted_number <- str_extract(n_helpful, pattern_2)
        review_n_helpful[k, ] <- as.numeric(extracted_number)
      }
    }

    k <- k + 1
  }

  # Use str_match on the vector review_dates
  dates <- str_match(review_dates, "on ([:alpha:]]+ [0-9]+, [0-9]+)")[, 2]

  # Convert the extracted dates to a standard date format
  formatted_dates <- as.Date(dates, format = "%B %d, %Y", locale = "en")

  # Convert the star ratings to numeric
  pattern_star <- "([0-9]+\\. [0-9]+) out of 5 stars"
  match_result_star <- str_match(star_ratings, pattern_star)[, 2]
  star_rating_num <- as.numeric(match_result_star)

  # Create a data frame with the extracted information
  reviews_comb <- data.frame(review_title, review_text, formatted_dates,
                             review_verified, review_n_helpful,
                             star_rating_num)

  # Append reviews to the data frame
  reviews_all <- rbind(reviews_all, reviews_comb)
} else {
  # Print a warning if the request fails
  warning(paste("Request failed with status code:",
                httr::status_code(response)))
}
reviews_scraped <- rbind(reviews_all)
# Add a random timeout to avoid being too aggressive
timeout <- runif(1, 5, 10)
Sys.sleep(timeout)
print(paste0(filter, " reviews: iteration ", i, "/", iters, " completed."))
}
}

save(reviews_scraped, file="Reviews_Scraped_500_v2.rda")

```

Now let us take a look at the dataframe that we have created via the code chunk above and the structure that we are working with:

```

load("reviews_scraped_500_v2.rda")
str(reviews_scraped)

```

```
## 'data.frame':    500 obs. of  6 variables:
```

```
## $ review_title    : chr  "Looks brand new and I love it! iPhone 6 to 11" "Surprisingly a go
## $ review_text     : chr  "My iPhone 6 died and was only iOS 10 so it was time to get a new p
## $ formatted_dates: Date, format: "2023-10-25" "2023-10-23" ...
## $ Verified       : logi  TRUE TRUE TRUE TRUE TRUE TRUE TRUE ...
## $ N_helpful      : num  9 19 3 3 3 2 1 5 1 NA ...
## $ star_rating_num: num  5 5 5 5 5 5 5 5 5 5 ...
```

Our dataframe consists of 500 observations (maximum number of reviews that can be scraped using our algorithm and not specific search words for Amazon) and 6 features. The features are the title of the review (character), the text of the review (character), the date the review was published (Date), whether it was a verified customer or not (TRUE / FALSE), the number of people that found it helpful and the number of stars the product was rated (both numeric).

In the following steps we will analyze these 500 reviews and their content.

General analysis

We want to start by taking a look at some general features of our review data to get a better understanding of the reviews in general.

Let us start by taking a look at **when the reviews were made**:

```
min(reviews_scraped$formatted_dates)
```

```
## [1] "2020-11-06"
```

```
max(reviews_scraped$formatted_dates)
```

```
## [1] "2023-11-17"
```

The time frame of our reviews spans from 06th of November, 2020, when the first review was published until the 17th of November, where the most recent one was made.

Additionally, let's check who made the reviewers, or to be exact: **was the reviewer a verified buyer** or not?

```
table(reviews_scraped$Verified)
```

```
##
## TRUE
## 500
```

This reveals that all 500 reviews were made by verified buyers, thus we can neglect this column.

What might be interesting is to have more of an insight into the **general rating behavior**. Since we have, already when scraping, turned the star rating into a numerical variable the following steps can be done quite easily:

```
min(reviews_scraped$star_rating_num)
```

```
## [1] 1
```

```
mean(reviews_scraped$star_rating_num)
```

```
## [1] 3
```

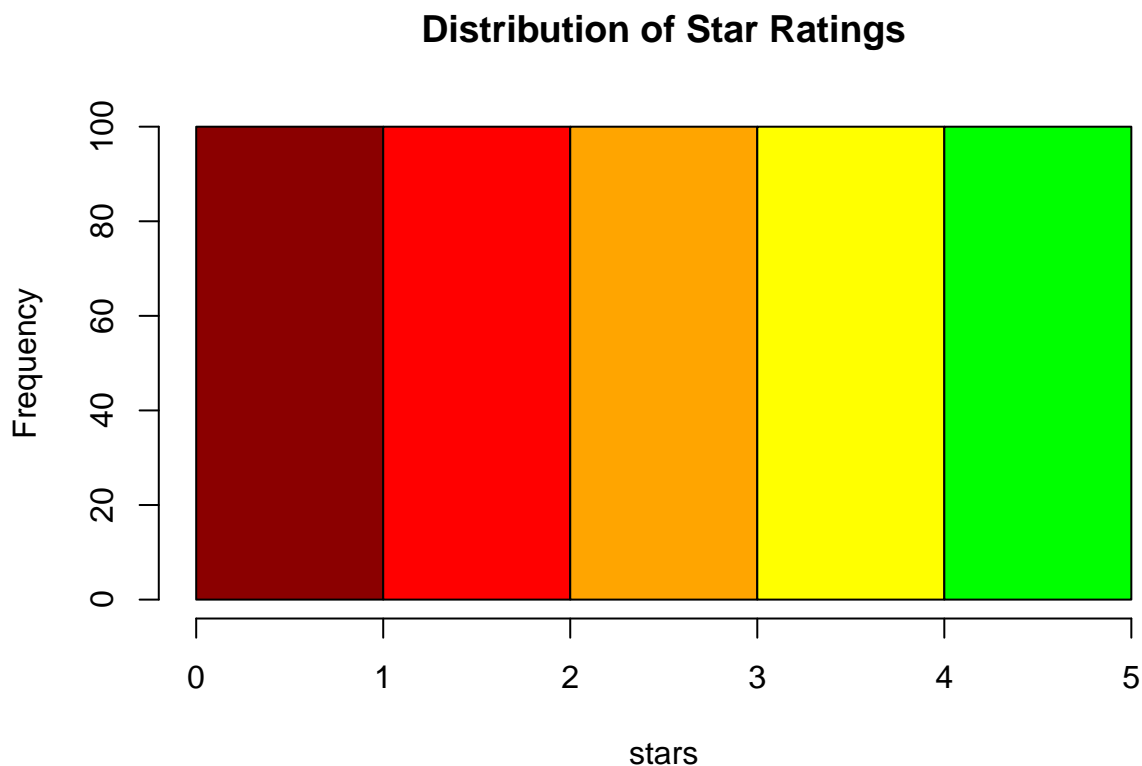
```
max(reviews_scraped$star_rating_num)
```

```
## [1] 5
```

The **worst rating** is 1 stars, the **highest one** is 5 stars and the **average rating** amounts to 3 stars. This is due to how we scraped the data.

Another visualization of the **distribution of star ratings** might be helpful to see how the product has performed in general. Thus, here we present a histogram:

```
stars <- reviews_scraped$star_rating_num
cols <- c("dark red", "red", "orange", "yellow", "green")
hist(stars, main="Distribution of Star Ratings", xlab="stars", col=cols,
     breaks = 0:5)
```



Here we can also see the bias introduced to the data due to the scraping method.

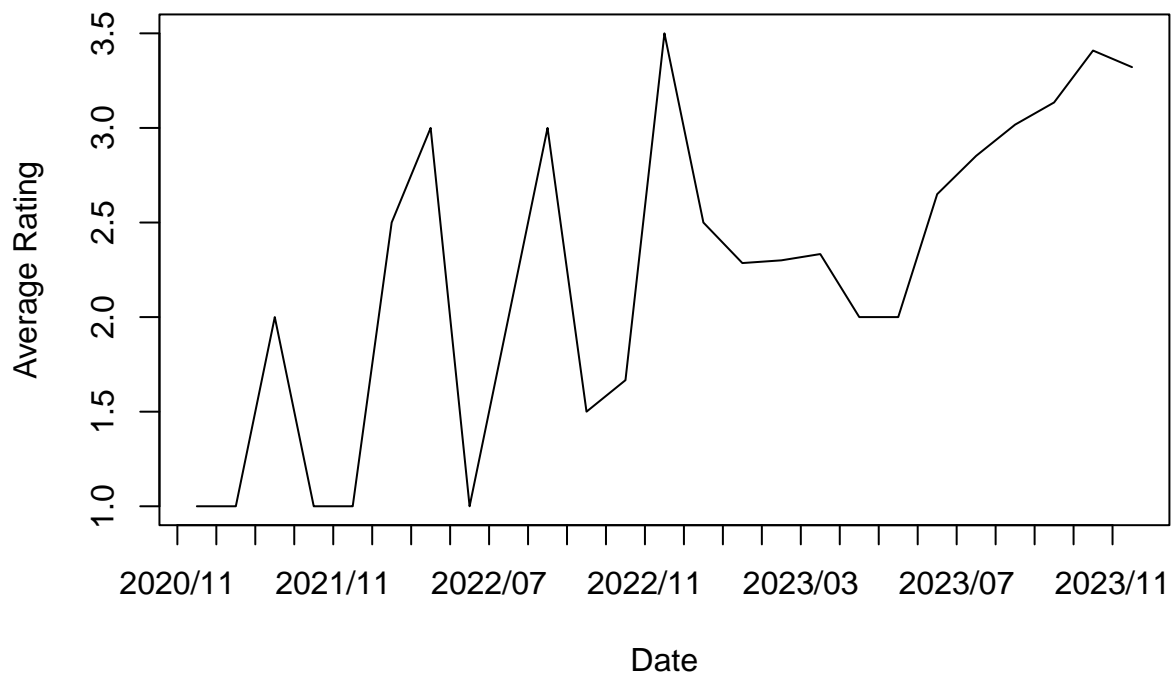
It might be also interesting to **reviews have changed over time**. For this we first summarize months.

```
dates <- strftime(reviews_scraped$formatted_dates, "%Y/%m")
```

Now we want to check whether the **average rating has improved or deteriorated** over time.

```
#plot rating distribution across time
plottingstars <- aggregate(reviews_scraped$star_rating_num ~ dates, FUN = mean)
plot(plottingstars[,2], type="l", xlab="Date", xaxt="n", ylab="Average Rating",
     ,main="Development of Average Star Ratings")
axis(side=1, at=1:nrow(plottingstars)-0.5,
     labels=plottingstars[1:nrow(plottingstars),1])
```

Development of Average Star Ratings

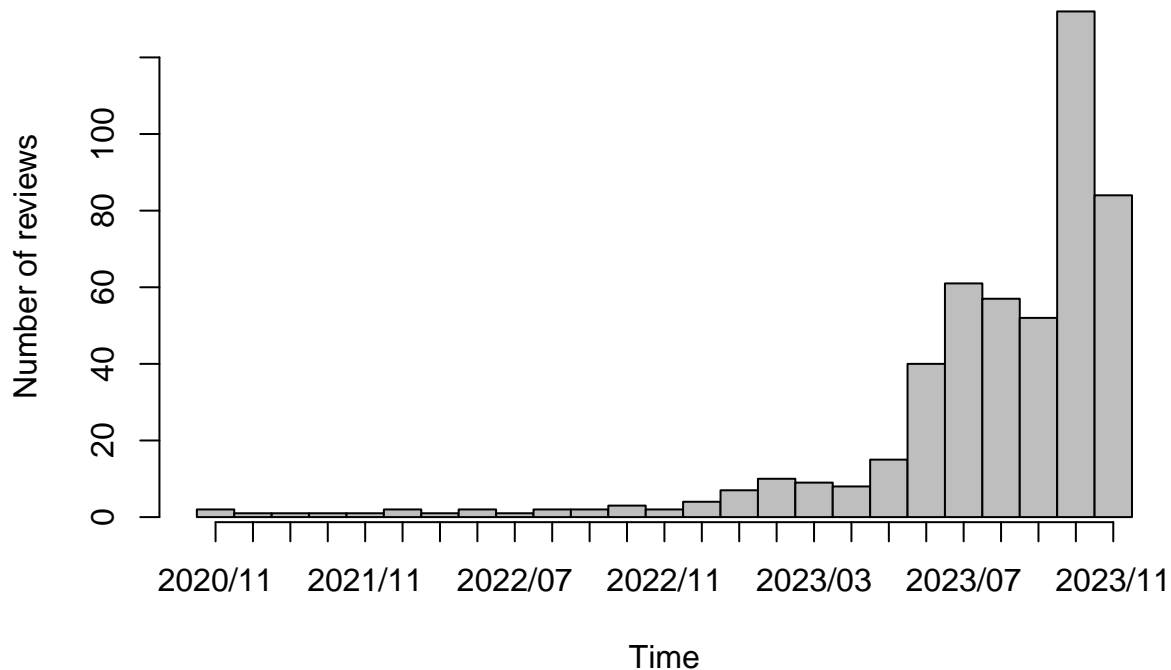


It seems as the average rating oscillated between Nov 2020 and Dec 2022, then shows a downwards trend until May 2023, after which it increases until the end of our time period.

It might also be interesting to take a look at **the number of reviews published across our time period** and see if and how this corresponds to the development of the star ratings.

```
# plot number of reviews distribution across time
plottingcount <- aggregate(reviews_scraped$star_rating_num ~ dates, FUN = length)
barplot(plottingcount[,2], main="Number of reviews across time", xlab="Time",
        xaxt="n", ylab="Number of reviews", space=0)
axis(side=1, at=1:nrow(plottingcount)-0.5,
     labels=plottingcount[1:nrow(plottingcount),1])
```

Number of reviews across time



As the plot suggests, our dataset contains mostly recent reviews.

Furthermore, we might want to take a look at **helpfulness of the reviews**. For this we start off with simply wanting to know how the reviews were distributed in this aspect.

```
reviews_scraped$N_helpful[is.na(reviews_scraped$N_helpful)] <- 0  
min(reviews_scraped$N_helpful)
```

```
## [1] 0
```

```
mean(reviews_scraped$N_helpful)
```

```
## [1] 2.088
```

```
median(reviews_scraped$N_helpful)
```

```
## [1] 0
```

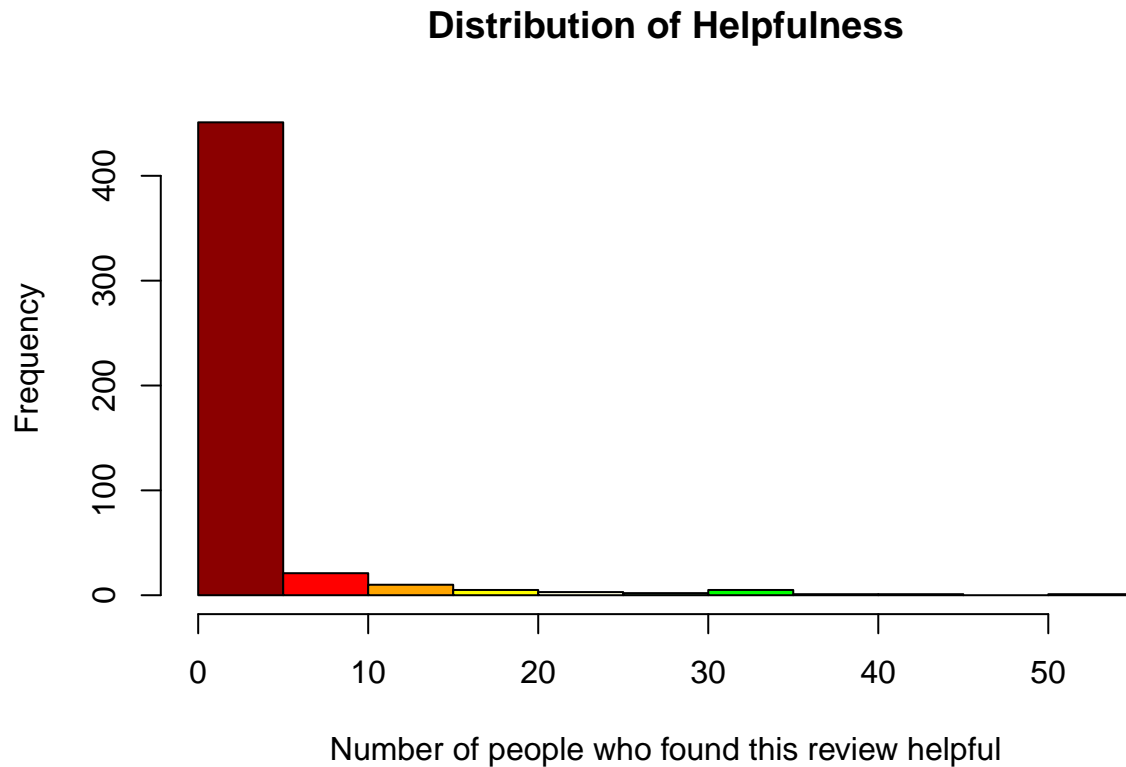
```
max(reviews_scraped$N_helpful)
```

```
## [1] 51
```

We can see that the reviews were rated as helpful by between 0 and 51 people, with an average of 2.088 and a median of 0 as well. This indicates skewedness and strong outliers, so we want to take a closer look.

```
helpful <- reviews_scraped$N_helpful  
cols <- c("dark red", "red", "orange", "yellow", "light yellow", "light green",  
          "green", "dark green", "blue", "purple")  
hist(helpful, main="Distribution of Helpfulness",
```

```
xlab="Number of people who found this review helpful", col=cols,
breaks = seq(0, 55, by=5))
```



This supports our hypothesis that the vast majority of reviews are rated as helpful by very few to no people. Only very few people attained over 15 “helpfulness”-votes.

This comes down to the following general helpfulness of reviews:

```
dim(reviews_scraped[reviews_scraped$N_helpful>0,])[1]/dim(reviews_scraped)[1]
```

```
## [1] 0.354
```

Overall a third of the reviews were perceived as helpful, meaning that at least 1 customer rated them as such.

Text Analysis

Now that we have taken a look at the ratings and developments over time of our reviews, let us dive deeper into **what the reviews actually say**:

In order to do some text analysis, we must first do some pre-processing of the reviews.

Let us start by looking at **how long the review texts are**:

```
reviews_scraped$review_text_2 <- reviews_scraped$review_text
reviews_scraped$review_text_2 <- gsub('[:punct:] '+' ',' ',
                                     reviews_scraped$review_text_2)
#split into substrings
reviews_scraped$review_text_2 <- strsplit(reviews_scraped$review_text_2, " ")
```

```
#count the number of strings = number of words
reviews_scraped$review_length <- sapply(reviews_scraped$review_text_2, length)
print(reviews_scraped$review_length)
```

```
## [1] 370 301 222 117 95 136 69 153 58 48 71 30 39 22 39 51 22 41
## [19] 40 35 49 23 20 50 53 78 19 273 5 14 13 67 48 24 18 48
## [37] 2 1 0 8 1 12 7 7 7 6 7 21 16 10 3 3 4 3
## [55] 14 65 14 59 50 11 17 11 4 23 84 2 16 52 25 34 3 10
## [73] 93 63 59 10 17 29 12 11 198 7 4 11 16 17 6 5 16 3
## [91] 3 17 11 2 2 42 22 13 1 27 422 197 78 192 339 39 57 76
## [109] 67 15 10 185 36 17 17 15 8 11 67 42 25 74 22 64 60 31
## [127] 12 26 2 12 84 63 20 16 10 45 53 126 55 3 29 11 15 52
## [145] 41 100 15 50 54 5 11 66 176 12 80 21 77 232 54 85 50 13
## [163] 21 22 8 5 12 15 56 100 28 9 27 131 31 62 16 4 166 2
## [181] 211 10 13 17 69 44 65 39 25 48 24 60 1 32 6 25 37 4
## [199] 5 103 220 390 169 50 37 88 32 79 27 26 24 76 23 24 66 26
## [217] 22 19 65 16 8 36 31 27 30 24 22 293 27 21 61 125 10 47
## [235] 83 12 102 37 5 54 92 4 61 47 4 84 248 26 181 45 17 95
## [253] 42 52 108 34 21 51 7 42 112 105 10 64 31 43 56 67 48 15
## [271] 40 42 175 25 12 31 145 86 65 140 132 116 65 32 112 67 1 34
## [289] 200 25 187 115 23 11 31 11 30 17 22 59 138 244 92 110 69 60
## [307] 55 285 53 52 43 64 64 209 56 33 136 30 48 29 44 113 19 55
## [325] 21 14 52 57 105 25 47 41 42 150 133 39 39 69 40 34 14 14
## [343] 80 80 53 91 52 77 66 65 184 21 82 42 22 40 82 29 17 174
## [361] 20 78 25 44 10 43 22 8 29 10 21 49 82 62 59 20 83 16
## [379] 137 49 168 59 7 27 21 134 65 19 50 9 86 18 72 183 77 37
## [397] 73 44 29 61 507 325 347 256 385 313 267 205 189 272 374 144 244 288
## [415] 133 115 164 108 202 227 126 323 183 266 82 84 201 159 80 76 94 171
## [433] 90 70 66 87 57 62 274 160 307 84 176 86 52 47 163 50 51 60
## [451] 47 45 43 42 41 39 45 168 42 38 212 35 35 72 192 32 66 204
## [469] 167 121 66 66 64 36 63 64 62 81 58 103 57 33 53 280 48 28
## [487] 31 174 96 210 111 46 48 105 29 27 28 86 96 27
```

Given this, we can for example now check the average length of reviews:

```
mean(reviews_scraped$review_length)
```

```
## [1] 70.392
```

We find out that the average review text is around 70 words long.

What might be more interesting is to see if there is some **correlation between the length of a reviews and its rating**:

```
correlation <- cor(reviews_scraped$review_length, reviews_scraped$star_rating_num)
correlation
```

```
## [1] -0.3355965
```

The correlation is -0.335. What this reveals is that long reviews tend to be more negative or that customers who have a negative opinion about a product tend to share more information to explain

or complain. This however is a weak correlation.

Let's see if there is some **correlation between the length of a review and its helpfulness**:

```
correlation <- cor(reviews_scraped$review_length, reviews_scraped$N_helpful)
correlation
```

```
## [1] 0.6302964
```

The correlation is at 0.63. This leads to the conclusion that longer reviews tend to be rated as more helpful as they for example share more information and details that can be useful for other potential customers.

Now let us dive even deeper into text analysis than merely the word count. Let's have a look at **what people actually write in their reviews**.

```
# we start by loading in our text data as a "corpus"
TextDoc <- Corpus(VectorSource((reviews_scraped$review_text_2)))

#Replacing "/", "@" and "/" with space
toSpace <- content_transformer(function(x, pattern) gsub(pattern, " ", x))
TextDoc <- tm_map(TextDoc, toSpace, "/")
```

```
## Warning in tm_map.SimpleCorpus(TextDoc, toSpace, "/"): transformation drops
## documents
```

```
TextDoc <- tm_map(TextDoc, toSpace, "@")
```

```
## Warning in tm_map.SimpleCorpus(TextDoc, toSpace, "@"): transformation drops
## documents
```

```
TextDoc <- tm_map(TextDoc, toSpace, "\\|")
```

```
## Warning in tm_map.SimpleCorpus(TextDoc, toSpace, "\\|"): transformation drops
## documents
```

```
# Convert the text to lower case
TextDoc <- tm_map(TextDoc, content_transformer(tolower))
```

```
## Warning in tm_map.SimpleCorpus(TextDoc, content_transformer(tolower)):
## transformation drops documents
```

```
# Remove numbers
TextDoc <- tm_map(TextDoc, removeNumbers)
```

```
## Warning in tm_map.SimpleCorpus(TextDoc, removeNumbers): transformation drops
## documents
```

```
# Eliminate extra white spaces
TextDoc <- tm_map(TextDoc, stripWhitespace)
```

```
## Warning in tm_map.SimpleCorpus(TextDoc, stripWhitespace): transformation drops
## documents
```



```

# Remove English common stopwords
TextDoc <- tm_map(TextDoc, removeWords, stopwords("english"))

## Warning in tm_map.SimpleCorpus(TextDoc, removeWords, stopwords("english")):
## transformation drops documents

# Text stemming - which reduces words to their root form
TextDoc <- tm_map(TextDoc, stemDocument)

## Warning in tm_map.SimpleCorpus(TextDoc, stemDocument): transformation drops
## documents

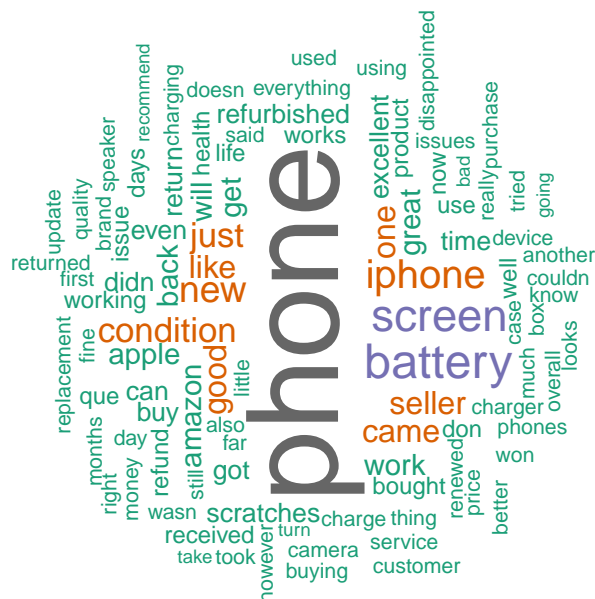
# Remove all punctuation
TextDoc <- tm_map(TextDoc, removePunctuation)

## Warning in tm_map.SimpleCorpus(TextDoc, removePunctuation): transformation
## drops documents

# we continue to now build a term-document matrix
TextDoc_tdm <- TermDocumentMatrix(TextDoc)
tdm_m <- as.matrix(TextDoc_tdm)
# we sort the terms by decreasing frequency
tdm_v <- sort(rowSums(tdm_m),decreasing=TRUE)
tdm_d <- data.frame(word = names(tdm_v),freq=tdm_v)

# with these steps accomplished we can now build a word cloud
set.seed(1234)
wordcloud(words = tdm_d$word, freq = tdm_d$freq, min.freq = 5,
          max.words=100, random.order=FALSE, rot.per=0.40,
          colors=brewer.pal(8, "Dark2"))

```



Words such as “iphone”, “phone” or “apple” do not surprise us as being frequently used since they refer to product name and brand. If we want to make room for new revelations, we can remove

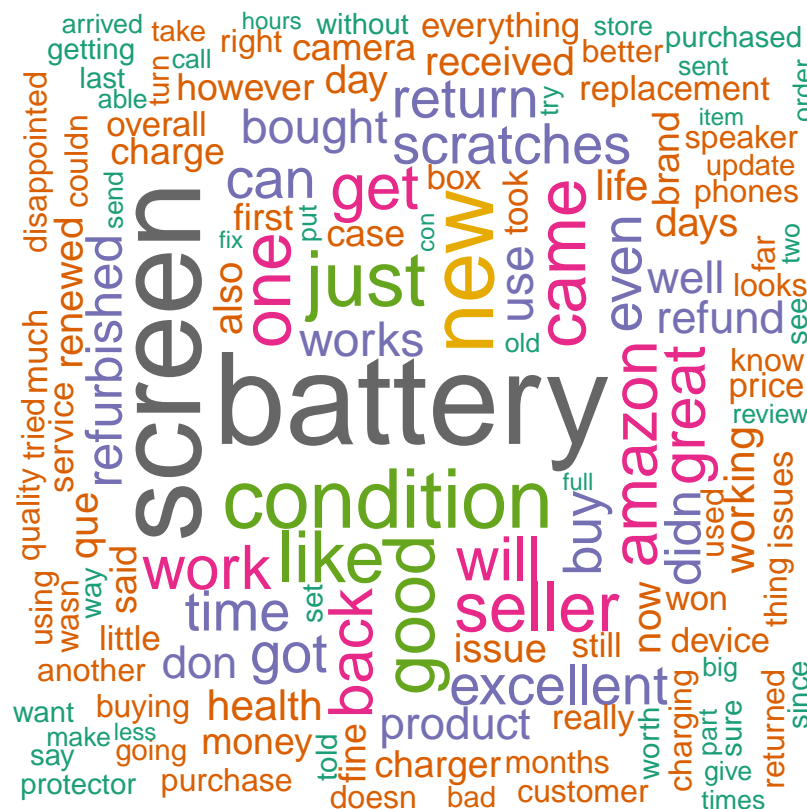
them.

Therefore we will make further alterations to remove chosen words such as the product name itself - this works via custom stop words:

```
# specify your custom stopwords as a character vector
TextDoc <- tm_map(TextDoc, removeWords, c("iphone", "phone", "apple"))

# Build a term-document matrix
TextDoc_tdm <- TermDocumentMatrix(TextDoc)
tdm_m <- as.matrix(TextDoc_tdm)
# Sort by decreasing value of frequency
tdm_v <- sort(rowSums(tdm_m), decreasing=TRUE)
tdm_d <- data.frame(word = names(tdm_v), freq=tdm_v)

#generate word cloud
set.seed(1234)
wordcloud(words = tdm_d$word, freq = tdm_d$freq, min.freq = 5,
          max.words=200, random.order=FALSE, rot.per=0.40,
          colors=brewer.pal(8, "Dark2"))
```



This reveals a much more interesting picture. We chose that all terms that are mentioned at least 5 times in the 500 reviews are shown in the word cloud. The by far biggest and thus most frequently used significant term in reviews is “battery”. Additionally screen and scratches as well as condition might be interesting ones to dive into.

Let us now do the same for the titles of the reviews and see if there are differences:

```

# we start by loading in our text data as a "corpus"
TextDoc <- Corpus(VectorSource((reviews_scraped$review_title)))

#Replacing "/", "@" and "|" with space
toSpace <- content_transformer(function (x , pattern ) gsub(pattern, " ", x))
TextDoc <- tm_map(TextDoc, toSpace, "/")
TextDoc <- tm_map(TextDoc, toSpace, "@")
TextDoc <- tm_map(TextDoc, toSpace, "\\|")
# Convert the text to lower case
TextDoc <- tm_map(TextDoc, content_transformer(tolower))
# Remove numbers
TextDoc <- tm_map(TextDoc, removeNumbers)
# Eliminate extra white spaces
TextDoc <- tm_map(TextDoc, stripWhitespace)
# Remove german common stopwords
TextDoc <- tm_map(TextDoc, removeWords, stopwords("english"))
# Text stemming - which reduces words to their root form
TextDoc <- tm_map(TextDoc, stemDocument)
# Remove all punctuation
TextDoc <- tm_map(TextDoc, removePunctuation)
# specify your custom stopwords as a character vector
TextDoc <- tm_map(TextDoc, removeWords, c("iphone","phone","apple"))

# Build a term-document matrix
TextDoc_tdm <- TermDocumentMatrix(TextDoc)
tdm_m <- as.matrix(TextDoc_tdm)
# Sort by decreasing value of frequency
tdm_v <- sort(rowSums(tdm_m),decreasing=TRUE)
tdm_d <- data.frame(word = names(tdm_v),freq=tdm_v)

# with these steps accomplished we can now build a word cloud
set.seed(1234)
wordcloud(words = tdm_d$word, freq = tdm_d$freq, min.freq = 2,
          max.words=100, random.order=FALSE, rot.per=0.40,
          colors=brewer.pal(8, "Dark2"))

```



Since this word cloud is less informative, we will include more words (at lower frequency threshold) and see an overwhelming frequency of very positive terms (good, grate, perfect). This is a positive indicator.

Sentiment analysis

Before we dive into applying sentiment analysis on our dataset of reviews, we tried to recall the basics of the Sentiment package and its meaning in R. Essentially sentiment analysis works by differentiating between words depending on the sentiment that is attached to them. This is conducted via dictionaries consisting of lists of positive vs. negative words, or lists of more diverse emotions. Packages such as `sentimentr` in R work by scanning the text to see if words in the text match with any dictionary entries. The words are then assigned a value (>0 if the word is located on the positive list, <0 if it is on the negative one) - all values are added together and the average sentiment is determined. Important note: The packages take the words before and after a term into account in order to assess its classification. This way valence shifters or negations can be included. Sentiment analysis can be applied in a number of fields and situations. Its use-cases range from social media monitoring, political campaigns, PR and market research.

We now want to start our first computations in the field of sentiment analysis to get a better picture about the general tonality and sentiment of the reviews we are examining. In this we will differentiate between title and text look at the correlation between them and their general behavior.

```
sentences <- get_sentences(reviews_scraped$review_text)
reviews_scraped$sentiment_text=sentiment_by(sentences)

sentences_title <- get_sentences(reviews_scraped$review_title)
reviews_scraped$sentiment_title=sentiment_by(sentences_title)

summary(cbind(reviews_scraped$sentiment_text$save_sentiment,
```

```

reviews_scraped$sentiment_title$ave_sentiment))

##           V1           V2
## Min.      :-0.54524   Min.      :-1.40729
## 1st Qu.: -0.04643   1st Qu.: -0.14434
## Median :  0.05303   Median :  0.00000
## Mean      :  0.11167   Mean      :  0.06231
## 3rd Qu.:  0.25000   3rd Qu.:  0.34820
## Max.      :  1.42500   Max.      :  1.23744

cor.test(reviews_scraped$sentiment_text$ave_sentiment,
         reviews_scraped$sentiment_title$ave_sentiment)

##
## Pearson's product-moment correlation
##
## data:  reviews_scraped$sentiment_text$ave_sentiment and reviews_scraped$sentiment_title$ave.
## t = 11.062, df = 498, p-value < 2.2e-16
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
##  0.3708724 0.5118749
## sample estimates:
##           cor
## 0.4441193

```

The computations reveal that in the review texts, sentiment ranges from -0.55 to 1.42 and averages at around 0.11. For the titles, the lowest sentiment is -1.4, the highest is 1.24 and the average is around 0.06. Both average sentiments are positive, which is good news for the product. Simply comparing the numbers allows us to make the inference that in general the sentiment in titles is more extreme than that in texts. This should not surprise us - titles are meant to catch people's attentions and thus, similar as headlines in the news, use more aggressive wordings, stronger opinions and more catchy phrases. We also prove that the sentiment of title and review text are correlated since the p-value is extremely low.

Now, similar to the way we wanted to check if the rating behavior changed over time we want to take a look at the **development of sentiment over time**:

```

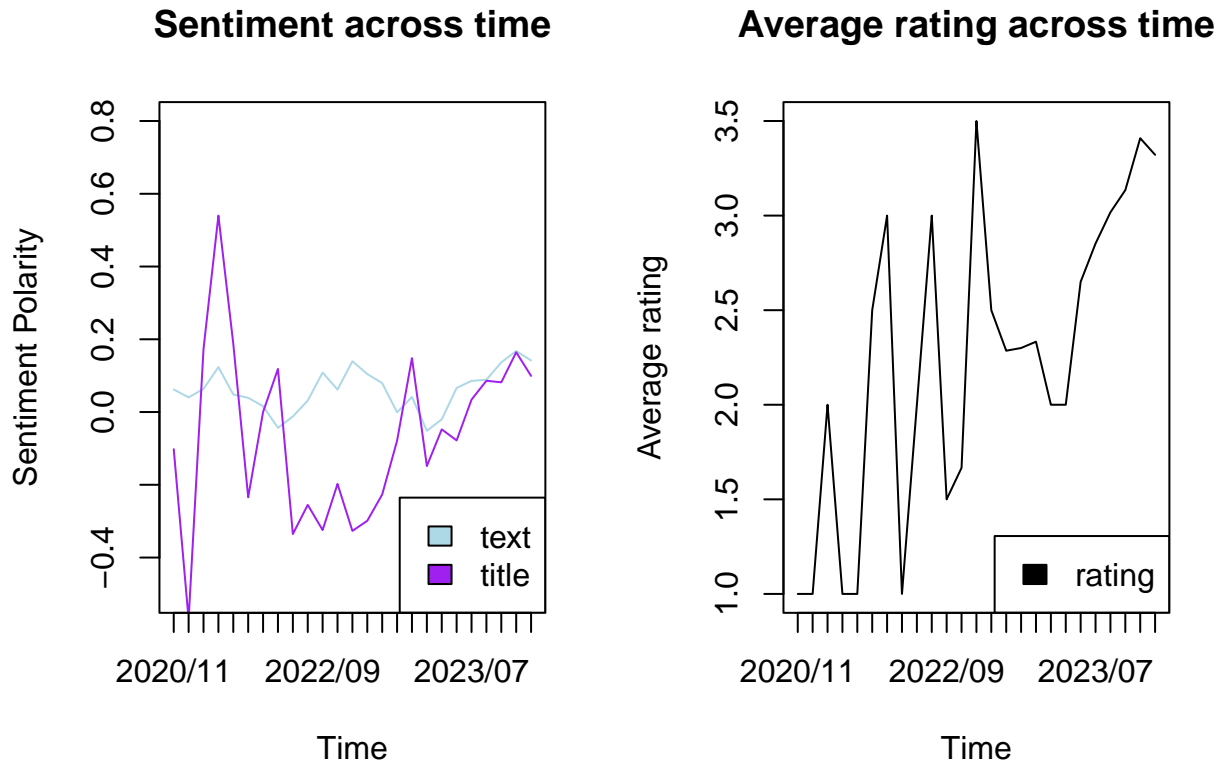
textsentiment_agg<-aggregate(reviews_scraped$sentiment_text$ave_sentiment ~ dates,
                             FUN = mean)
titlesentiment_agg<-aggregate(reviews_scraped$sentiment_title$ave_sentiment ~ dates,
                              FUN = mean)
starratings_agg<-aggregate(reviews_scraped$star_rating_num ~ dates, FUN = mean)

par(mfrow=c(1,2))
plot(textsentiment_agg[,2],type="l",xlab="Time",xaxt="n",
     ylab="Sentiment Polarity", ylim=c(-0.5, 0.8),
     main="Sentiment across time",col="light blue")
lines(titlesentiment_agg[,2],type="l",xlab="Time",xaxt="n",
     ylab="Sentiment Polarity", main="Sentiment across time",col="purple")

```

```
axis(side=1,at=1:nrow(textsentiment_agg),
      labels=textsentiment_agg[1:nrow(textsentiment_agg),1])
legend("bottomright",c("text","title"),fill=c("light blue","purple"))

plot(starratings_agg[,2],type="l",xlab="Time",xaxt="n",ylab="Average rating",
      main="Average rating across time")
axis(side=1,at=1:nrow(starratings_agg),
      labels=starratings_agg[1:nrow(starratings_agg),1])
legend("bottomright","rating",fill="black")
```



When comparing the plots of how sentiment (in both title and text) developed, we can see the same trends between title sentiment and ratings. Additionally the left plot highlights the extremity of titles compared to text.

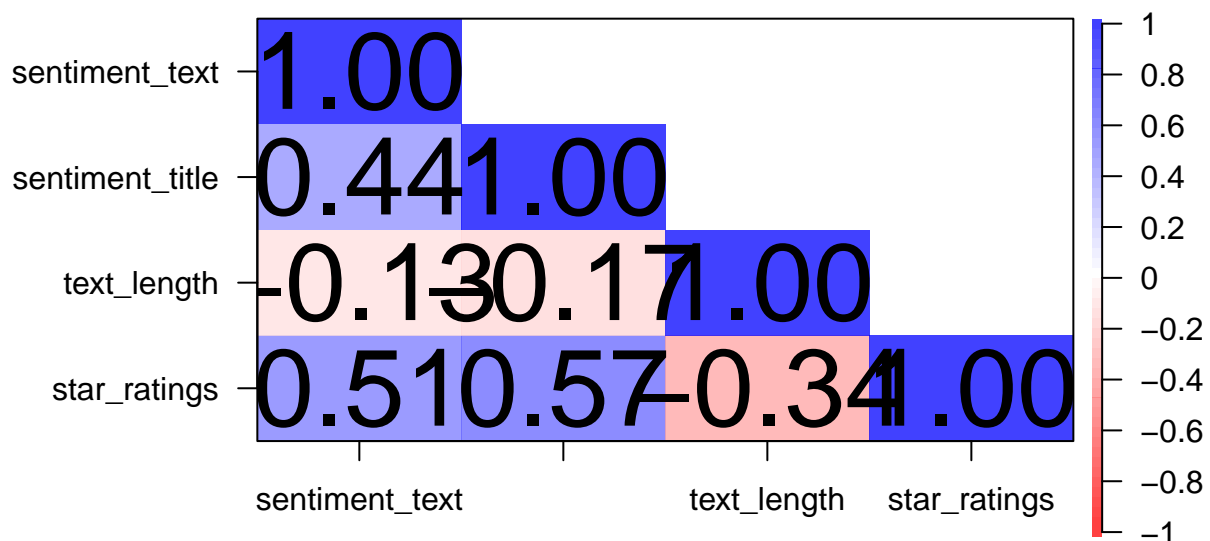
We can take a look at how sentiment & other variables interact in more detail:

```
corrm <- cor(cbind(reviews_scraped$sentiment_text$ave_sentiment,
                   reviews_scraped$sentiment_title$ave_sentiment,
                   reviews_scraped$review_length, reviews_scraped$star_rating_num))
colnames(corrm) <- c("sentiment_text","sentiment_title","text_length",
                    "star_ratings")
rownames(corrm) <- c("sentiment_text","sentiment_title","text_length",
                    "star_ratings")
corPlot(corrm ,numbers=T,diag=T,upper=F, main="Correlations between variables",
        xact = "n")
```

```
## Warning in axis(2, at = at2, labels = lab2, las = ylas, ...): "xact" is not a
```

```
## graphical parameter
## Warning in axis(xaxis, at = at1, labels = lab1, las = xlas, line = line, :
## "xact" is not a graphical parameter
## Warning in text.default(rx, ry, rv, cex = 1.5 * cex, ...): "xact" is not a
## graphical parameter
## Warning in axis(4, at = at2, labels = labels, las = 2, ...): "xact" is not a
## graphical parameter
```

Correlations between variables



From the correlation plot we can see a rather strong positive correlation between sentiment in title and sentiment in text - this we have already found out. Additionally we can see slight negative correlations between text length and sentiment. This further supports our findings that longer reviews will have a worse rating. Star ratings thus negatively correlate with text length, but are positively impacted by the sentiment score.

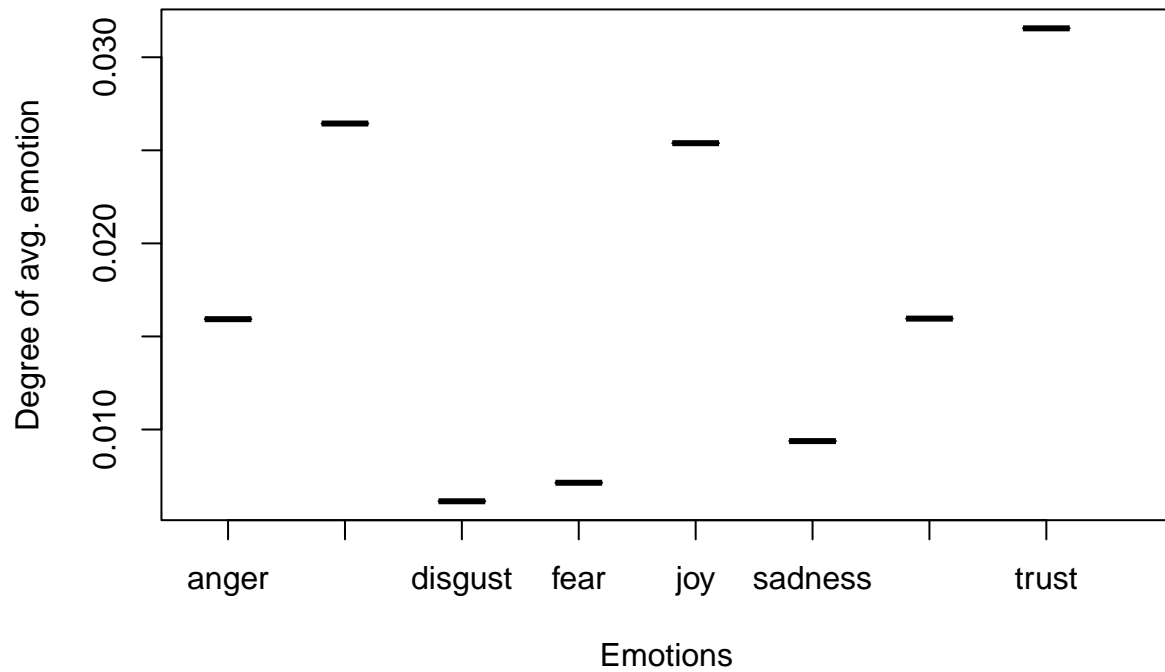
Emotion Analysis

In emotion analysis we get more detailed insights in the emotions that are expressed in reviews by differentiating not only between positively and negatively annotated terms, but a more developed spectrum of emotions. As a basis, Plutchik's wheel of emotion is used.

We now want to find out about the emotions in our reviews:

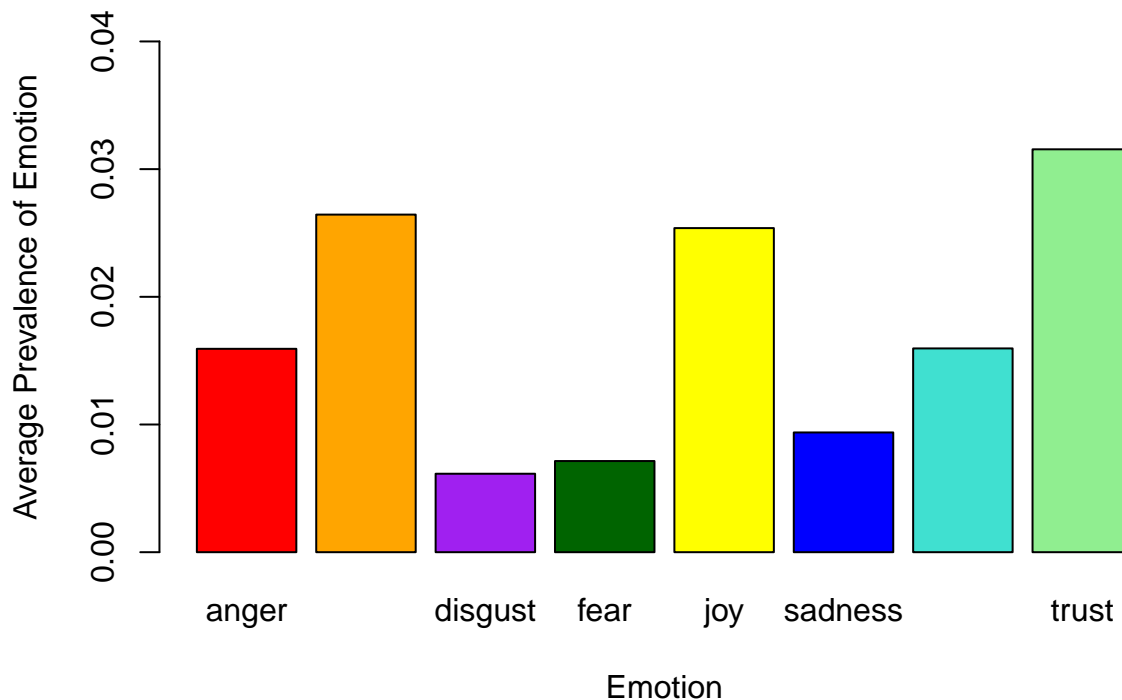
```
text_emotions <- emotion_by(sentences)
emotiontext_agg<-aggregate(text_emotions$aave_emotion ~ (text_emotions$emotion_type),
                           FUN = mean)
emotiontext_agg <- subset(emotiontext_agg,
                          (1:dim(emotiontext_agg)[1])%%2==1) #non-negated only
plot(emotiontext_agg,xaxt="n", ylab="Degree of avg. emotion",xlab="Emotions",
     main="Wheel of Emotions for Review Texts")
axis(side=1,at=row.names(emotiontext_agg),labels=emotiontext_agg[c(1:8),1])
```

Wheel of Emotions for Review Texts



```
colors = c("red", "orange", "purple", "dark green", "yellow", "blue", "turquoise",  
           "light green")  
barplot(emotiontext_agg[, 2], names.arg = emotiontext_agg[, 1], col = colors,  
        main="Emotion Analysis for Review Texts", ylim=c(0, 0.04), xlab="Emotion",  
        ylab = "Average Prevalence of Emotion")
```


Emotion Analysis for Review Texts



Essentially, this plot reveals information about how dominant certain emotions are in our review text. The higher the bar of a certain emotion, the higher its relevance. In the colored bar plot visualizations we can see which emotions correspond to which part of the wheel and can infer that the most dominant emotions in the review texts are anticipation (orange), joy (yellow) and trust (light green).

Now let us take a closer look at the emotions separately:

```
#add emotions to data set
#reduce to non-negated emotions
text_emotions_nnegated <- text_emotions[-grep("negated",
                                              text_emotions$emotion_type),]

#get average emotions for each review
temp <- reshape(text_emotions_nnegated[,c(1,2,6)], idvar = "element_id",
               timevar = "emotion_type", v.names = "ave_emotion",
               direction = "wide")
reviews_scraped <- cbind(reviews_scraped,temp)
emotions_list <- c("ave_emotion.anger", "ave_emotion.anticipation",
                  "ave_emotion.disgust", "ave_emotion.fear", "ave_emotion.joy",
                  "ave_emotion.sadness", "ave_emotion.surprise",
                  "ave_emotion.trust")
summary(reviews_scraped[,emotions_list])
```

```
## ave_emotion.anger ave_emotion.anticipation ave_emotion.disgust
## Min. :0.000000 Min. :0.00000 Min. :0.000000
## 1st Qu.:0.000000 1st Qu.:0.00000 1st Qu.:0.000000
```

```
## Median :0.001193 Median :0.01299 Median :0.000000
## Mean :0.015928 Mean :0.02644 Mean :0.006145
## 3rd Qu.:0.022990 3rd Qu.:0.03071 3rd Qu.:0.006501
## Max. :0.200000 Max. :1.00000 Max. :0.093750
## ave_emotion.fear ave_emotion.joy ave_emotion.sadness ave_emotion.surprise
## Min. :0.000000 Min. :0.000000 Min. :0.000000 Min. :0.000000
## 1st Qu.:0.000000 1st Qu.:0.000000 1st Qu.:0.000000 1st Qu.:0.000000
## Median :0.000000 Median :0.009788 Median :0.000000 Median :0.000000
## Mean :0.007139 Mean :0.025383 Mean :0.009378 Mean :0.01596
## 3rd Qu.:0.010583 3rd Qu.:0.029155 3rd Qu.:0.015152 3rd Qu.:0.01515
## Max. :0.111111 Max. :1.000000 Max. :0.084746 Max. :1.00000
## ave_emotion.trust
## Min. :0.00000
## 1st Qu.:0.00000
## Median :0.01852
## Mean :0.03155
## 3rd Qu.:0.03704
## Max. :1.00000
```

This again shows that trust, anticipation and joy have the highest means, thus are on average the prevalent emotions expressed in the reviews we are looking at.

Now we want to check if and how emotional categories are potentially related to one another:

```
corrE <- cor(reviews_scraped[,emotions_list])
colnames(corrE) <- c("anger","anticipation","disgust","fear","joy","sadness",
                    "surprise","trust")
rownames(corrE) <-c("anger","anticipation","disgust","fear","joy","sadness",
                  "surprise","trust")

corPlot(corrE ,numbers=T,diag=F,upper=F, main="Correlations between emotions",
        xact="n" )
```

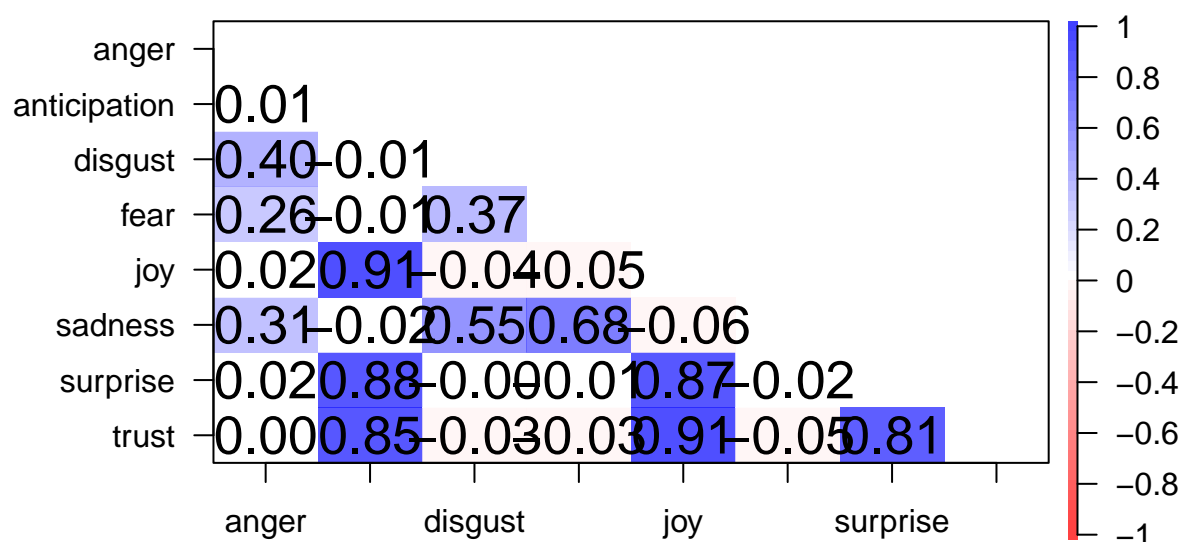
```
## Warning in axis(2, at = at2, labels = lab2, las = ylas, ...): "xact" is not a
## graphical parameter
```

```
## Warning in axis(xaxis, at = at1, labels = lab1, las = xlas, line = line, :
## "xact" is not a graphical parameter
```

```
## Warning in text.default(rx, ry, rv, cex = 1.5 * cex, ...): "xact" is not a
## graphical parameter
```

```
## Warning in axis(4, at = at2, labels = labels, las = 2, ...): "xact" is not a
## graphical parameter
```

Correlations between emotions



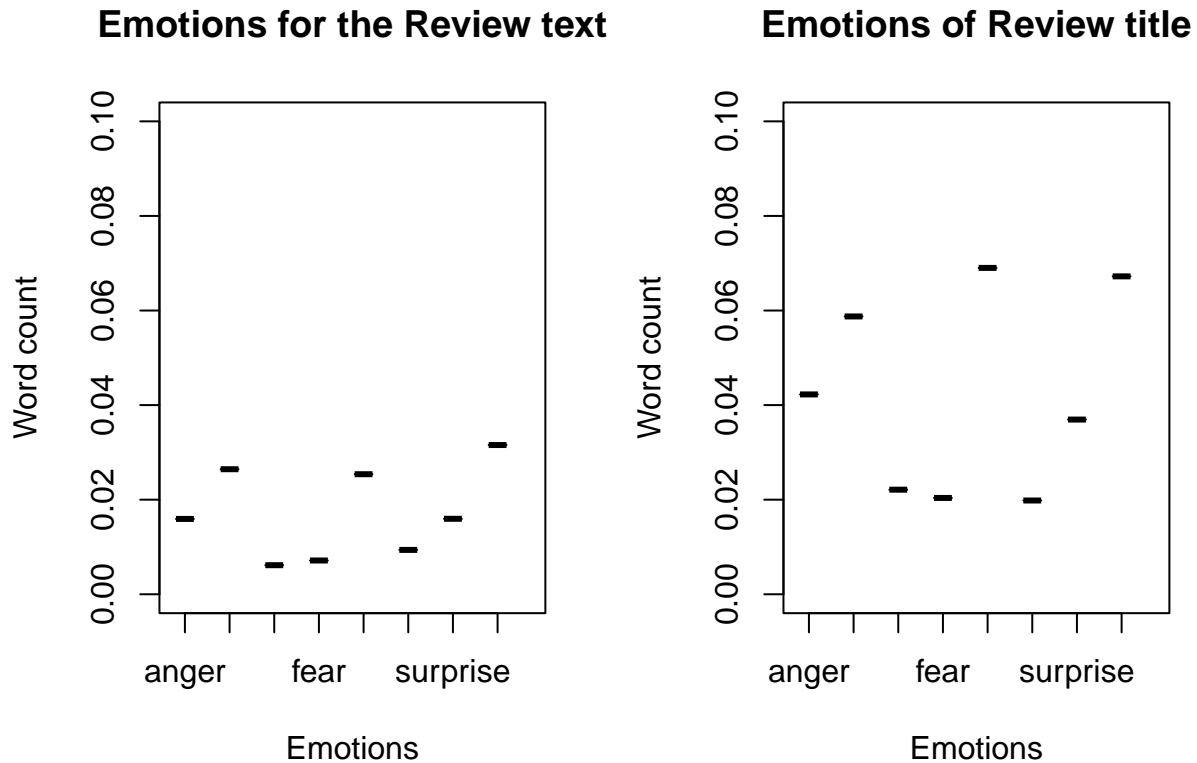
The correlation matrix provides us with information about how the different emotions might go hand in hand with another in the reviews we are examining. This reveals a positive correlation of anger with disgust, fear and sadness; of anticipation with joy, surprise and trust; and of joy with trust and surprise. Sadness also correlated with anger disgust and fear. Negative correlations are not significant.

Let us now also check if our emotions are similar between text and title:

```
title_emotions <- emotion_by(sentences_title)
emotiontitle_agg<-aggregate(title_emotions$ave_emotion ~ (title_emotions$emotion_type),
                             FUN = mean)

emotiontitle_agg <- subset(emotiontitle_agg,(1:dim(emotiontitle_agg)[1])%%2==1)

par(mfrow=c(1,2))
plot(emotiontext_agg,xaxt="n",ylab="Word count",xlab="Emotions",
     main="Emotions for the Review text",
     ylim=c(0,0.1))
axis(side=1,at=row.names(emotiontext_agg),labels=emotiontext_agg[c(1:8),1])
plot(emotiontitle_agg,xaxt="n",ylab="Word count",xlab="Emotions",
     main="Emotions of Review title",
     ylim=c(0,0.1))
axis(side=1,at=row.names(emotiontitle_agg),labels=emotiontitle_agg[c(1:8),1])
```



We can see that for anticipation, joy, and trust are the main three emotions of the title, with all emotions being higher than for the text. This shows how much more emotional the titles are.

Topic Analysis

Now that we have talked about the sentiments and emotions prevalent in the reviews of the product, we want to take a look at the content of the reviews. To do this efficiently, we will make use of the methods of Topic Analysis.

Topic Analysis refers to methods that aim to identify the different contents discussed and help us focus on those that we are really interested in. It functions via LDA (Latent Dirichlet allocation): after pre-processing, a document term matrix is created, the number of topics to differentiate between is decided upon, we check for convergence, estimate the model & can then interpret the first conclusion.

```
# already pre-processed in steps for creating the wordcloud
DocText_dtm <- DocumentTermMatrix(TextDoc)
```

```
# first steps
raw.sum<- apply(DocText_dtm,1,FUN=sum)
DocText_dtm <- DocText_dtm[raw.sum!=0,]
DocText_dtm
```

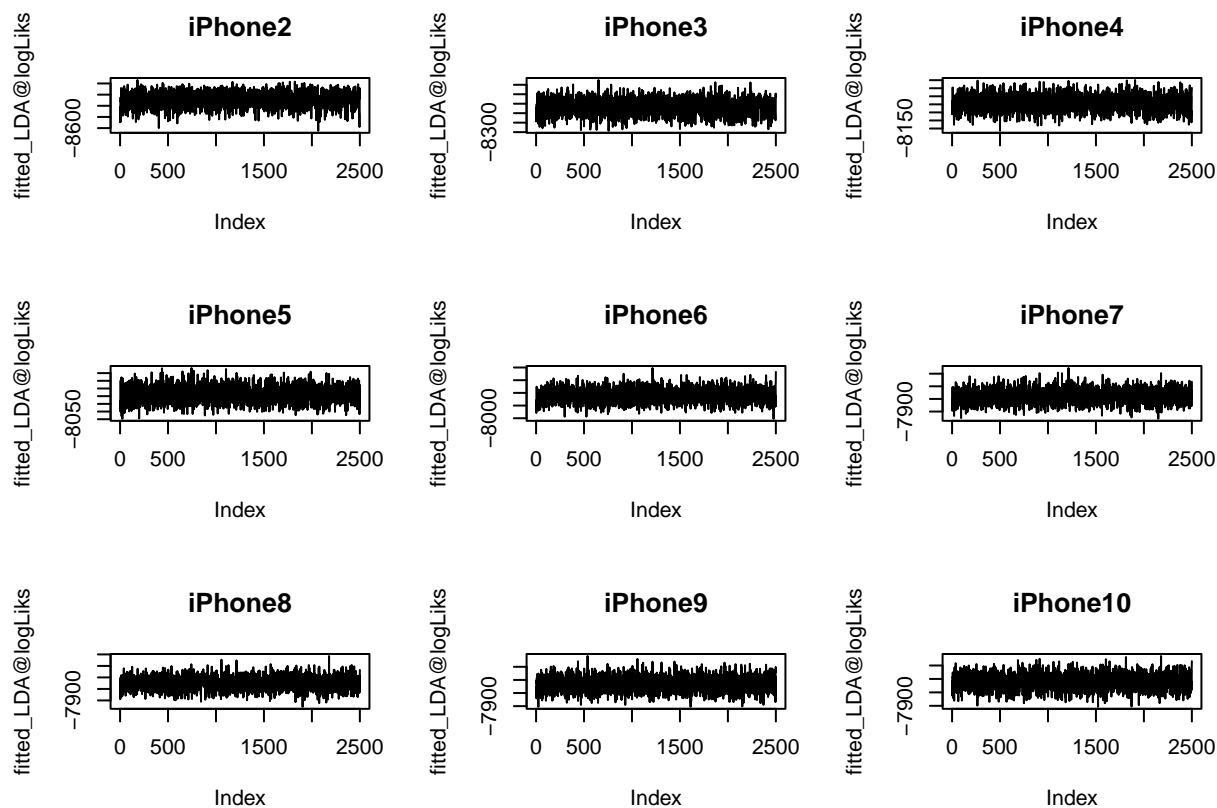
```
## <<DocumentTermMatrix (documents: 488, terms: 577)>>
## Non-/sparse entries: 1398/280178
## Sparsity           : 100%
## Maximal term length: 16
## Weighting          : term frequency (tf)
```

One of the trickiest part of topic analysis is deciding upon the right number of topics to distinguish between. To compute the optimal number, we will run the following code:

```
SEED <- 123
burnin <- 5000 #not being used for estimation
iter <- 20000 #number of iterations after burn in
keep <- 10 #keeps every tenth iteration (=burin + iter)
maxtops <- 10
avg_result_fin <- matrix(nrow=maxtops-1,ncol=3)
counter <- 1
par(mfrow=c(3,3))
for (k in 2:maxtops){

  fitted_LDA <- LDA(DocText_dtm, k = k, method = "Gibbs",
                    control = list(seed = SEED,burnin = burnin, iter = iter,
                                   keep = keep) )
  plot(fitted_LDA@logLiks,type="l", main=paste0(c("iPhone",k),collapse=""))
  words_LDA <- dim(posterior(fitted_LDA)[[1]])[2]
  avg_result_fin[counter,] <- cbind(k, logLik(fitted_LDA),
                                    -2*logLik(fitted_LDA)+(k+k*words_LDA))

  counter=counter+1
}
```



This for loop iterates and finds the best number of topics for our dataset.

```
colnames(avg_result_fin) <- c("ntopics","ll","AIC")
avg_result_fin
```

```
##      ntopics      ll      AIC
## [1,]      2 -8421.724 17999.45
## [2,]      3 -8105.016 17944.03
## [3,]      4 -7897.937 18107.87
## [4,]      5 -7826.560 18543.12
## [5,]      6 -7634.700 18737.40
## [6,]      7 -7663.824 19373.65
## [7,]      8 -7635.237 19894.47
## [8,]      9 -7640.158 20482.32
## [9,]     10 -7568.215 20916.43
```

We make our decision on the number of topics based on AIC. We want to minimize AIC and thus choose 3 topics.

As the next step, we estimate our topic model.

```
ktop <- 3

fitted_LDA_model <- LDA(DocText_dtm, k = ktop, method = "Gibbs",
                        control = list(seed = SEED, burnin = burnin,
                                       iter = iter, keep = keep) )
```

From our model, we extract the topics.

```
topics <- posterior(fitted_LDA_model)[[1]]
```

Afterwards we look at the ten most important words for each topic in order to get a grasp on what the topics are about.

```
#look at the ten most important words for the topics
for (k in 1:ktop){
  print(sort(topics[k,],decreasing=T)[1:10])
}
```

```
##      good      work      iphon      like      excel      's      purchas
## 0.09529907 0.08249497 0.05688678 0.02944942 0.02396195 0.01847448 0.01298701
##      issu      condition      excelent
## 0.01298701 0.01115786 0.01115786
##      batteri      buy      perfect      qualiti      speaker      iphon      use
## 0.07307465 0.03762064 0.02383297 0.01989364 0.01989364 0.01792397 0.01595430
##      scratch      defect      nice
## 0.01595430 0.01595430 0.01398464
##      great      screen      new      condit      life      worth      seller
## 0.07070707 0.04211931 0.04021346 0.03259005 0.02496665 0.01734324 0.01734324
##      price      replac      get
## 0.01353154 0.01162569 0.01162569
```

We label topic 1 as **quality**, topic 2 as **hardware** due to the high proportion of comments mentioning

hardware aspects and potential issues, and topic 3 as **value**, since the product we are analysing is refurbished, and the reviews mention words like condition, worth, seller and price.

Next we create the dataframe multi, which shows the relative importance of the three topics for each review.

```
multi <- posterior(fitted_LDA_model)[[2]]  
dim(multi)
```

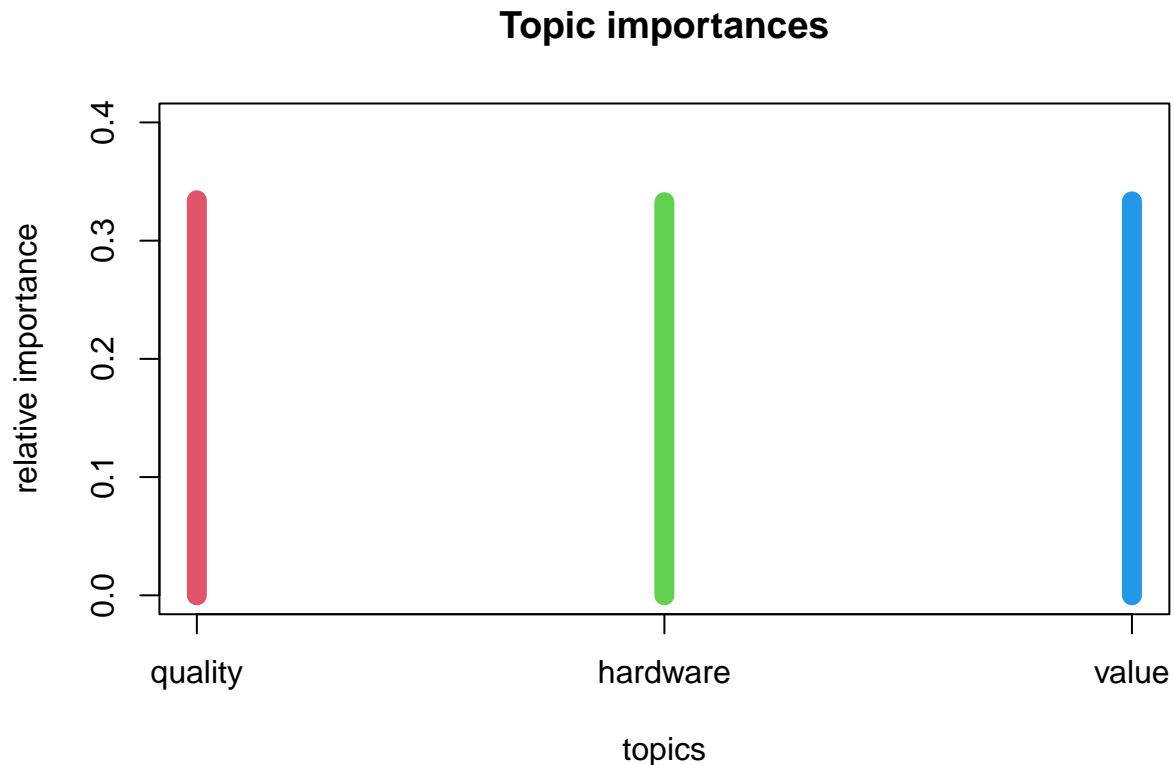
```
## [1] 488    3
```

```
colnames(multi) <- c("quality", "hardware", "value")  
summary(multi)
```

```
##      quality      hardware      value  
##  Min.   :0.2825   Min.   :0.2924   Min.   :0.2924  
## 1st Qu.:0.3205   1st Qu.:0.3205   1st Qu.:0.3205  
## Median :0.3333   Median :0.3272   Median :0.3306  
## Mean   :0.3342   Mean   :0.3326   Mean   :0.3333  
## 3rd Qu.:0.3457   3rd Qu.:0.3397   3rd Qu.:0.3397  
## Max.   :0.3827   Max.   :0.3908   Max.   :0.3869
```

We plot the average importance of the topics.

```
plot(colSums(multi)/dim(multi)[1],type="h", main="Topic importances",ylim=c(0,.4),  
     ylab="relative importance",xlab="topics",xaxt="n",lwd=10,col=2:4)  
axis(side=1,at=1:3,labels=colnames(multi))
```



From the plot above we can observe that the relative importance per review of each topic is about

the same (a third), meaning the three topics are of equal importance in the reviews. Looking back on the summary statistics above, we can also see that the minimum importances of all topics are just under 0.3, and the maxima for all three are below 0.4, which implies that all of our reviews include all three topics relatively evenly.

It is worth to note however, that since we are working with a limited amount of data (488 reviews), it can easily be the case that we simply do not have enough data to clearly identify topics.

Modeling approaches

In the following we will attempt to model the ratings of our product using different features that we have derived in the previous chapters.

Modelling ratings based on **date**:

```
stars_1 <- lm(star_rating_num ~ as.factor(strftime(formatted_dates, "%m")),
              data=reviews_scraped)
summary(stars_1)
```

```
##
## Call:
## lm(formula = star_rating_num ~ as.factor(strftime(formatted_dates,
##           "%m")), data = reviews_scraped)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.353  -1.017   0.000   1.175   2.000
##
## Coefficients:
##                                Estimate Std. Error t value
## (Intercept)                   2.28571     0.52142   4.384
## as.factor(strftime(formatted_dates, "%m"))02  0.01429     0.67985   0.021
## as.factor(strftime(formatted_dates, "%m"))03  0.07792     0.66701   0.117
## as.factor(strftime(formatted_dates, "%m"))04 -0.28571     0.71399  -0.400
## as.factor(strftime(formatted_dates, "%m"))05 -0.28571     0.61954  -0.461
## as.factor(strftime(formatted_dates, "%m"))06  0.28571     0.56320   0.507
## as.factor(strftime(formatted_dates, "%m"))07  0.53968     0.54963   0.982
## as.factor(strftime(formatted_dates, "%m"))08  0.73123     0.55149   1.326
## as.factor(strftime(formatted_dates, "%m"))09  0.78836     0.55419   1.423
## as.factor(strftime(formatted_dates, "%m"))10  1.06723     0.53467   1.996
## as.factor(strftime(formatted_dates, "%m"))11  0.96148     0.54154   1.775
## as.factor(strftime(formatted_dates, "%m"))12  0.21429     0.86468   0.248
##                                Pr(>|t|)
## (Intercept)                   1.43e-05 ***
## as.factor(strftime(formatted_dates, "%m"))02  0.9832
## as.factor(strftime(formatted_dates, "%m"))03  0.9070
## as.factor(strftime(formatted_dates, "%m"))04  0.6892
## as.factor(strftime(formatted_dates, "%m"))05  0.6449
## as.factor(strftime(formatted_dates, "%m"))06  0.6122
## as.factor(strftime(formatted_dates, "%m"))07  0.3266
```



```
## as.factor(strftime(formatted_dates, "%m"))08 0.1855
## as.factor(strftime(formatted_dates, "%m"))09 0.1555
## as.factor(strftime(formatted_dates, "%m"))10 0.0465 *
## as.factor(strftime(formatted_dates, "%m"))11 0.0764 .
## as.factor(strftime(formatted_dates, "%m"))12 0.8044
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.38 on 488 degrees of freedom
## Multiple R-squared:  0.07125,    Adjusted R-squared:  0.05032
## F-statistic: 3.404 on 11 and 488 DF,  p-value: 0.000145
```

This model takes a look at the star ratings as a function of time. We use the months as factor variables with 02 - February representing our baseline. However, we find that only October has a sufficiently low p-value to be considered significant.

Modelling ratings based on **review length**:

```
stars_2 <- lm(star_rating_num ~ review_length, data=reviews_scraped)
summary(stars_2)

##
## Call:
## lm(formula = star_rating_num ~ review_length, data = reviews_scraped)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.2663 -1.1481 -0.1435  0.9823  3.8386
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   3.4319674  0.0807214   42.52 < 2e-16 ***
## review_length -0.0061366  0.0007719   -7.95 1.26e-14 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.335 on 498 degrees of freedom
## Multiple R-squared:  0.1126, Adjusted R-squared:  0.1108
## F-statistic: 63.21 on 1 and 498 DF,  p-value: 1.257e-14
```

This model takes a look at how the review length impacts the star ratings in closer detail. We find out that actually, the review length does have a significant impact in determining the stars given in a review. This relationship is negative, meaning: the longer a review, the lower the average star rating.

Modelling ratings based on **sentiment**:

```
stars_3 <- lm((star_rating_num) ~ (sentiment_title$ave_sentiment)
+ (sentiment_text$ave_sentiment)
+ (review_length),
data=reviews_scraped)
```

```
summary(stars_3)
```

```
##
## Call:
## lm(formula = (star_rating_num) ~ (sentiment_title$ave_sentiment) +
##      (sentiment_text$ave_sentiment) + (review_length), data = reviews_scraped)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.62866 -0.79767  0.02889  0.79394  2.85312
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      3.0229077   0.0682189   44.312 < 2e-16 ***
## sentiment_title$ave_sentiment  1.3429874   0.1253157   10.717 < 2e-16 ***
## sentiment_text$ave_sentiment   1.6706981   0.2006733    8.325 8.23e-16 ***
## review_length    -0.0041647   0.0006127   -6.797 3.08e-11 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.042 on 496 degrees of freedom
## Multiple R-squared:  0.4613, Adjusted R-squared:  0.458
## F-statistic: 141.6 on 3 and 496 DF,  p-value: < 2.2e-16
```

This model now includes the sentiment of the title of the review and of the review itself as well as the length of the review as determining parameters for the expected star rating. In the case of this model, the sentiment of title as well as the review length are considered to be significant. While a higher sentiment of the title will lead to an increase in the predicted star rating, a higher word count of the review will lead to a decrease in the predicted star rating. We also tested for a model only using sentiment, not review length - here too, only the sentiment of the title resulted in being statistically significant at $\alpha=0.05$.

Modelling ratings based on **emotions**:

```
stars_4 <- lm((star_rating_num ~ ave_emotion.anger+ave_emotion.anticipation +
              ave_emotion.disgust +ave_emotion.fear +ave_emotion.joy +
              ave_emotion.sadness+ave_emotion.surprise+ave_emotion.trust +
              (review_length)), data=reviews_scraped)
summary(stars_4)
```

```
##
## Call:
## lm(formula = (star_rating_num ~ ave_emotion.anger + ave_emotion.anticipation +
##      ave_emotion.disgust + ave_emotion.fear + ave_emotion.joy +
##      ave_emotion.sadness + ave_emotion.surprise + ave_emotion.trust +
##      (review_length)), data = reviews_scraped)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
```

```
## -2.6921 -1.0501 -0.1458 0.9583 3.5966
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    3.436e+00  9.518e-02  36.097 < 2e-16 ***
## ave_emotion.anger    1.799e+00  2.394e+00   0.751  0.45287
## ave_emotion.anticipation -3.487e+00  2.246e+00  -1.553  0.12110
## ave_emotion.disgust   -1.412e+01  5.279e+00  -2.675  0.00772 **
## ave_emotion.fear     -5.898e+00  5.395e+00  -1.093  0.27483
## ave_emotion.joy       8.135e+00  2.758e+00   2.950  0.00333 **
## ave_emotion.sadness   -6.303e+00  5.643e+00  -1.117  0.26460
## ave_emotion.surprise  -3.240e+00  2.161e+00  -1.499  0.13448
## ave_emotion.trust     1.135e+00  2.017e+00   0.563  0.57365
## review_length     -5.323e-03  7.525e-04  -7.074 5.23e-12 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.276 on 490 degrees of freedom
## Multiple R-squared:  0.2018, Adjusted R-squared:  0.1871
## F-statistic: 13.76 on 9 and 490 DF,  p-value: < 2.2e-16
```

This model reveals that stars are significantly impacted by review length, and only two of the emotions: joy and disgust.

Modelling ratings based on **sentiments & emotion**:

```
stars_5 <- lm((star_rating_num ~ ave_emotion.anger+ave_emotion.anticipation +
              ave_emotion.disgust +ave_emotion.fear +ave_emotion.joy +
              ave_emotion.sadness+ave_emotion.surprise+ave_emotion.trust +
              (sentiment_title$ave_sentiment) +
              (sentiment_text$ave_sentiment) +
              review_length),
              data=reviews_scraped)
summary(stars_5)
```

```
##
## Call:
## lm(formula = (star_rating_num ~ ave_emotion.anger + ave_emotion.anticipation +
##      ave_emotion.disgust + ave_emotion.fear + ave_emotion.joy +
##      ave_emotion.sadness + ave_emotion.surprise + ave_emotion.trust +
##      (sentiment_title$ave_sentiment) + (sentiment_text$ave_sentiment) +
##      review_length), data = reviews_scraped)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.75800 -0.77881  0.01105  0.75331  2.89580
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
```

```
## (Intercept)                3.0638350  0.0827222  37.038 < 2e-16 ***
## ave_emotion.anger          1.0415091  1.9524080   0.533   0.594
## ave_emotion.anticipation   -2.8119720  1.8365760  -1.531   0.126
## ave_emotion.disgust        -0.5169613  4.3975422  -0.118   0.906
## ave_emotion.fear           -4.0902856  4.4024474  -0.929   0.353
## ave_emotion.joy            2.5449571  2.2767434   1.118   0.264
## ave_emotion.sadness        -2.9618250  4.6260288  -0.640   0.522
## ave_emotion.surprise       -1.3880474  1.7661439  -0.786   0.432
## ave_emotion.trust          1.1092350  1.6438798   0.675   0.500
## sentiment_title$ave_sentiment 1.3190432  0.1262619  10.447 < 2e-16 ***
## sentiment_text$ave_sentiment 1.5839373  0.2234456   7.089 4.77e-12 ***
## review_length              -0.0040125  0.0006195  -6.477 2.28e-10 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.04 on 488 degrees of freedom
## Multiple R-squared:  0.4717, Adjusted R-squared:  0.4598
## F-statistic: 39.61 on 11 and 488 DF,  p-value: < 2.2e-16
```

When modelling star ratings based on emotions, sentiments & review_length we find out that no emotions significant impact in determining the final rating, only the sentiment scores and word count of the review.

```
stars_5a <- lm((star_rating_num ~ ave_emotion.anger+ave_emotion.anticipation +
  ave_emotion.disgust +ave_emotion.fear +ave_emotion.joy +
  ave_emotion.sadness+ave_emotion.surprise+ave_emotion.trust +
  (sentiment_title$ave_sentiment) + (sentiment_text$ave_sentiment) +
  review_length + as.factor(strftime(formatted_dates, "%m"))),
  data=reviews_scraped)
summary(stars_5a)
```

```
##
## Call:
## lm(formula = (star_rating_num ~ ave_emotion.anger + ave_emotion.anticipation +
##   ave_emotion.disgust + ave_emotion.fear + ave_emotion.joy +
##   ave_emotion.sadness + ave_emotion.surprise + ave_emotion.trust +
##   (sentiment_title$ave_sentiment) + (sentiment_text$ave_sentiment) +
##   review_length + as.factor(strftime(formatted_dates, "%m"))),
##   data = reviews_scraped)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.80656 -0.80960  0.03766  0.77116  2.81469
##
## Coefficients:
##                                Estimate Std. Error t value
## (Intercept)                   3.1599275   0.4142832   7.627
## ave_emotion.anger              0.9768932   1.9772068   0.494
## ave_emotion.anticipation       -2.6656199   1.8568833  -1.436
```

```

## ave_emotion.disgust          -0.1455391  4.4621573  -0.033
## ave_emotion.fear             -3.7363495  4.4743110  -0.835
## ave_emotion.joy              2.3222831  2.3069693   1.007
## ave_emotion.sadness          -3.1638402  4.6857719  -0.675
## ave_emotion.surprise         -1.2780611  1.7900549  -0.714
## ave_emotion.trust            1.0694200  1.6583713   0.645
## sentiment_title$ave_sentiment 1.3183986  0.1284598  10.263
## sentiment_text$ave_sentiment  1.5377999  0.2266045   6.786
## review_length               -0.0038090  0.0006597  -5.774
## as.factor(strftime(formatted_dates, "%m"))02 -0.1028738  0.5171225  -0.199
## as.factor(strftime(formatted_dates, "%m"))03 -0.4084537  0.5087765  -0.803
## as.factor(strftime(formatted_dates, "%m"))04 -0.2626961  0.5457888  -0.481
## as.factor(strftime(formatted_dates, "%m"))05 -0.4771553  0.4725151  -1.010
## as.factor(strftime(formatted_dates, "%m"))06 -0.1918187  0.4305455  -0.446
## as.factor(strftime(formatted_dates, "%m"))07 -0.2170265  0.4235811  -0.512
## as.factor(strftime(formatted_dates, "%m"))08 -0.0627838  0.4250594  -0.148
## as.factor(strftime(formatted_dates, "%m"))09 -0.0722265  0.4264306  -0.169
## as.factor(strftime(formatted_dates, "%m"))10 -0.0346479  0.4136791  -0.084
## as.factor(strftime(formatted_dates, "%m"))11 -0.0451920  0.4194122  -0.108
## as.factor(strftime(formatted_dates, "%m"))12  0.2228051  0.6586534   0.338
##                               Pr(>|t|)
## (Intercept)                  1.31e-13 ***
## ave_emotion.anger             0.621
## ave_emotion.anticipation      0.152
## ave_emotion.disgust           0.974
## ave_emotion.fear              0.404
## ave_emotion.joy               0.315
## ave_emotion.sadness           0.500
## ave_emotion.surprise          0.476
## ave_emotion.trust             0.519
## sentiment_title$ave_sentiment < 2e-16 ***
## sentiment_text$ave_sentiment  3.42e-11 ***
## review_length                 1.39e-08 ***
## as.factor(strftime(formatted_dates, "%m"))02  0.842
## as.factor(strftime(formatted_dates, "%m"))03  0.422
## as.factor(strftime(formatted_dates, "%m"))04  0.631
## as.factor(strftime(formatted_dates, "%m"))05  0.313
## as.factor(strftime(formatted_dates, "%m"))06  0.656
## as.factor(strftime(formatted_dates, "%m"))07  0.609
## as.factor(strftime(formatted_dates, "%m"))08  0.883
## as.factor(strftime(formatted_dates, "%m"))09  0.866
## as.factor(strftime(formatted_dates, "%m"))10  0.933
## as.factor(strftime(formatted_dates, "%m"))11  0.914
## as.factor(strftime(formatted_dates, "%m"))12  0.735
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.046 on 477 degrees of freedom

```

```
## Multiple R-squared:  0.4777, Adjusted R-squared:  0.4536
## F-statistic: 19.83 on 22 and 477 DF,  p-value: < 2.2e-16
```

With the inclusion of dates the significant predictors remain unchanged: sentiment scores and review length.

```
stars_step <- step(stars_5a)
```

```
## Start:  AIC=67.8
## star_rating_num ~ ave_emotion.anger + ave_emotion.anticipation +
##   ave_emotion.disgust + ave_emotion.fear + ave_emotion.joy +
##   ave_emotion.sadness + ave_emotion.surprise + ave_emotion.trust +
##   (sentiment_title$ave_sentiment) + (sentiment_text$ave_sentiment) +
##   review_length + as.factor(strftime(formatted_dates, "%m"))
##
```

	Df	Sum of Sq	RSS	AIC
## - as.factor(strftime(formatted_dates, "%m"))	11	6.013	528.30	51.526
## - ave_emotion.disgust	1	0.001	522.29	65.804
## - ave_emotion.anger	1	0.267	522.55	66.059
## - ave_emotion.trust	1	0.455	522.74	66.239
## - ave_emotion.sadness	1	0.499	522.78	66.281
## - ave_emotion.surprise	1	0.558	522.84	66.337
## - ave_emotion.fear	1	0.764	523.05	66.533
## - ave_emotion.joy	1	1.110	523.39	66.864
## <none>			522.29	67.803
## - ave_emotion.anticipation	1	2.256	524.54	67.958
## - review_length	1	36.507	558.79	99.585
## - sentiment_text\$ave_sentiment	1	50.426	572.71	111.886
## - sentiment_title\$ave_sentiment	1	115.332	637.62	165.565

```
## Step:  AIC=51.53
## star_rating_num ~ ave_emotion.anger + ave_emotion.anticipation +
##   ave_emotion.disgust + ave_emotion.fear + ave_emotion.joy +
##   ave_emotion.sadness + ave_emotion.surprise + ave_emotion.trust +
##   sentiment_title$ave_sentiment + sentiment_text$ave_sentiment +
##   review_length
##
```

	Df	Sum of Sq	RSS	AIC
## - ave_emotion.disgust	1	0.015	528.31	49.540
## - ave_emotion.anger	1	0.308	528.61	49.818
## - ave_emotion.sadness	1	0.444	528.74	49.946
## - ave_emotion.trust	1	0.493	528.79	49.992
## - ave_emotion.surprise	1	0.669	528.97	50.159
## - ave_emotion.fear	1	0.934	529.23	50.410
## - ave_emotion.joy	1	1.353	529.65	50.805
## <none>			528.30	51.526
## - ave_emotion.anticipation	1	2.538	530.84	51.922
## - review_length	1	45.418	573.72	90.763
## - sentiment_text\$ave_sentiment	1	54.399	582.70	98.530

```

## - sentiment_title$ave_sentiment 1 118.150 646.45 150.442
##
## Step: AIC=49.54
## star_rating_num ~ ave_emotion.anger + ave_emotion.anticipation +
## ave_emotion.fear + ave_emotion.joy + ave_emotion.sadness +
## ave_emotion.surprise + ave_emotion.trust + sentiment_title$ave_sentiment +
## sentiment_text$ave_sentiment + review_length
##
##
## Df Sum of Sq RSS AIC
## - ave_emotion.anger 1 0.295 528.61 47.819
## - ave_emotion.trust 1 0.490 528.80 48.004
## - ave_emotion.sadness 1 0.586 528.90 48.094
## - ave_emotion.surprise 1 0.671 528.98 48.175
## - ave_emotion.fear 1 0.928 529.24 48.418
## - ave_emotion.joy 1 1.360 529.67 48.826
## <none> 528.31 49.540
## - ave_emotion.anticipation 1 2.552 530.86 49.949
## - review_length 1 45.404 573.72 88.764
## - sentiment_text$ave_sentiment 1 55.805 584.12 97.747
## - sentiment_title$ave_sentiment 1 119.195 647.51 149.262
##
## Step: AIC=47.82
## star_rating_num ~ ave_emotion.anticipation + ave_emotion.fear +
## ave_emotion.joy + ave_emotion.sadness + ave_emotion.surprise +
## ave_emotion.trust + sentiment_title$ave_sentiment + sentiment_text$ave_sentiment +
## review_length
##
##
## Df Sum of Sq RSS AIC
## - ave_emotion.sadness 1 0.454 529.06 46.248
## - ave_emotion.trust 1 0.468 529.08 46.261
## - ave_emotion.surprise 1 0.673 529.28 46.455
## - ave_emotion.fear 1 0.861 529.47 46.633
## - ave_emotion.joy 1 1.433 530.04 47.173
## <none> 528.61 47.819
## - ave_emotion.anticipation 1 2.596 531.20 48.268
## - review_length 1 46.525 575.13 87.996
## - sentiment_text$ave_sentiment 1 55.639 584.25 95.857
## - sentiment_title$ave_sentiment 1 119.058 647.67 147.383
##
## Step: AIC=46.25
## star_rating_num ~ ave_emotion.anticipation + ave_emotion.fear +
## ave_emotion.joy + ave_emotion.surprise + ave_emotion.trust +
## sentiment_title$ave_sentiment + sentiment_text$ave_sentiment +
## review_length
##
##
## Df Sum of Sq RSS AIC
## - ave_emotion.trust 1 0.492 529.55 44.713
## - ave_emotion.surprise 1 0.678 529.74 44.888

```

```

## - ave_emotion.joy          1      1.407 530.47  45.576
## <none>                      529.06  46.248
## - ave_emotion.anticipation  1      2.640 531.70  46.737
## - ave_emotion.fear         1      3.268 532.33  47.327
## - review_length            1     46.624 575.69  86.477
## - sentiment_text$ave_sentiment 1     59.396 588.46  97.448
## - sentiment_title$ave_sentiment 1    118.939 648.00 145.642
##
## Step:  AIC=44.71
## star_rating_num ~ ave_emotion.anticipation + ave_emotion.fear +
##   ave_emotion.joy + ave_emotion.surprise + sentiment_title$ave_sentiment +
##   sentiment_text$ave_sentiment + review_length
##
##              Df Sum of Sq    RSS    AIC
## - ave_emotion.surprise      1      0.617 530.17  43.295
## <none>                      529.55  44.713
## - ave_emotion.anticipation    1      2.421 531.97  44.994
## - ave_emotion.fear            1      3.212 532.77  45.736
## - ave_emotion.joy            1      3.756 533.31  46.246
## - review_length              1     46.601 576.15  84.884
## - sentiment_text$ave_sentiment 1     59.447 589.00  95.909
## - sentiment_title$ave_sentiment 1    118.923 648.48 144.009
##
## Step:  AIC=43.29
## star_rating_num ~ ave_emotion.anticipation + ave_emotion.fear +
##   ave_emotion.joy + sentiment_title$ave_sentiment + sentiment_text$ave_sentiment +
##   review_length
##
##              Df Sum of Sq    RSS    AIC
## <none>                      530.17  43.295
## - ave_emotion.joy           1      3.152 533.32  44.258
## - ave_emotion.fear           1      3.248 533.42  44.349
## - ave_emotion.anticipation    1      4.408 534.58  45.434
## - review_length              1     46.376 576.55  83.223
## - sentiment_text$ave_sentiment 1     60.194 590.36  95.066
## - sentiment_title$ave_sentiment 1    119.584 649.75 142.993
summary(stars_step)

##
## Call:
## lm(formula = star_rating_num ~ ave_emotion.anticipation + ave_emotion.fear +
##   ave_emotion.joy + sentiment_title$ave_sentiment + sentiment_text$ave_sentiment +
##   review_length, data = reviews_scraped)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.65434 -0.78976 -0.00045  0.74003  2.81941

```



```
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    3.0740814   0.0740868  41.493 < 2e-16 ***
## ave_emotion.anticipation -3.3199554   1.6398799  -2.025  0.0435 *
## ave_emotion.fear      -5.7536598   3.3107169  -1.738  0.0829 .
## ave_emotion.joy       2.9765486   1.7386884   1.712  0.0875 .
## sentiment_title$ave_sentiment 1.3216012   0.1253278  10.545 < 2e-16 ***
## sentiment_text$ave_sentiment 1.6194492   0.2164585   7.482 3.40e-13 ***
## review_length    -0.0040343   0.0006143  -6.567 1.31e-10 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.037 on 493 degrees of freedom
## Multiple R-squared:  0.4698, Adjusted R-squared:  0.4634
## F-statistic: 72.82 on 6 and 493 DF,  p-value: < 2.2e-16
```

We pick our best model using AIC:

```
AIC(stars_1)
```

```
## [1] 1754.552
```

```
AIC(stars_2)
```

```
## [1] 1711.768
```

```
AIC(stars_3)
```

```
## [1] 1466.226
```

```
AIC(stars_4)
```

```
## [1] 1674.837
```

```
AIC(stars_5)
```

```
## [1] 1472.465
```

```
AIC(stars_5a)
```

```
## [1] 1488.741
```

```
AIC(stars_step)
```

```
## [1] 1464.233
```

stars_3 has the lowest score outside the one with stepwise selection, which only uses sentiments and review length. stars_step includes anticipation, fear and joy as well, further lowering the AIC by a slight amount.

We do model for determining if the rating will have 5 stars:

```
stars_6 <- glm((star_rating_num) > 4 ~ (ave_emotion.anger+ave_emotion.anticipation +
                                          ave_emotion.disgust +ave_emotion.fear +
```

```

ave_emotion.joy + ave_emotion.sadness +
ave_emotion.surprise + ave_emotion.trust +
sentiment_title$ave_sentiment +
sentiment_text$ave_sentiment+review_length),
data=reviews_scraped, family = "binomial")
summary(stars_6)

##
## Call:
## glm(formula = (star_rating_num) > 4 ~ (ave_emotion.anger + ave_emotion.anticipation +
##     ave_emotion.disgust + ave_emotion.fear + ave_emotion.joy +
##     ave_emotion.sadness + ave_emotion.surprise + ave_emotion.trust +
##     sentiment_title$ave_sentiment + sentiment_text$ave_sentiment +
##     review_length), family = "binomial", data = reviews_scraped)
##
## Coefficients:
##
##             Estimate Std. Error z value Pr(>|z|)
## (Intercept)    -1.829541    0.270671  -6.759 1.39e-11 ***
## ave_emotion.anger    -10.764945    6.621895  -1.626  0.1040
## ave_emotion.anticipation    -4.088539    5.248829  -0.779  0.4360
## ave_emotion.disgust    -7.646071   18.940134  -0.404  0.6864
## ave_emotion.fear   -30.611894   20.771997  -1.474  0.1406
## ave_emotion.joy     10.428342    6.727616   1.550  0.1211
## ave_emotion.sadness     7.279685   17.188276   0.424  0.6719
## ave_emotion.surprise   -14.172392    7.091098  -1.999  0.0456 *
## ave_emotion.trust     5.336966    3.703566   1.441  0.1496
## sentiment_title$ave_sentiment    2.850839    0.425870   6.694 2.17e-11 ***
## sentiment_text$ave_sentiment    1.928125    0.595520   3.238  0.0012 **
## review_length    -0.006096    0.002657  -2.294  0.0218 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##     Null deviance: 500.40  on 499  degrees of freedom
## Residual deviance: 333.02  on 488  degrees of freedom
## AIC: 357.02
##
## Number of Fisher Scoring iterations: 6

```

When combining all features in a model, besides the above mentioned predictors, surprise is also significant, negatively impacting the odds of a five star review. No other emotions are significant.

```

stars_6a <- glm((star_rating_num) > 4 ~ (ave_emotion.anger+ave_emotion.anticipation +
ave_emotion.disgust +ave_emotion.fear +
ave_emotion.joy + ave_emotion.sadness +
ave_emotion.surprise+ave_emotion.trust +
sentiment_title$ave_sentiment +

```

```

                                sentiment_text$ave_sentiment+review_length +
                                as.factor(strftime(formatted_dates, "%m"))),
                                data=reviews_scraped, family = "binomial")
summary(stars_6a)

##
## Call:
## glm(formula = (star_rating_num) > 4 ~ (ave_emotion.anger + ave_emotion.anticipation +
##     ave_emotion.disgust + ave_emotion.fear + ave_emotion.joy +
##     ave_emotion.sadness + ave_emotion.surprise + ave_emotion.trust +
##     sentiment_title$ave_sentiment + sentiment_text$ave_sentiment +
##     review_length + as.factor(strftime(formatted_dates, "%m"))),
##     family = "binomial", data = reviews_scraped)
##
## Coefficients:
##
##               Estimate Std. Error z value
## (Intercept)      -2.008e+01  5.904e+03  -0.003
## ave_emotion.anger      -1.031e+01  8.096e+00  -1.274
## ave_emotion.anticipation -2.334e+00  8.150e+00  -0.286
## ave_emotion.disgust       1.085e+01  2.064e+01   0.526
## ave_emotion.fear       -1.321e+01  2.505e+01  -0.527
## ave_emotion.joy         1.397e+01  8.859e+00   1.576
## ave_emotion.sadness     -7.588e+00  2.117e+01  -0.358
## ave_emotion.surprise    -1.826e+01  9.900e+00  -1.845
## ave_emotion.trust        4.321e+00  4.020e+00   1.075
## sentiment_title$ave_sentiment  3.830e+00  5.999e-01   6.385
## sentiment_text$ave_sentiment  1.769e+00  7.407e-01   2.388
## review_length      -2.443e-03  2.621e-03  -0.932
## as.factor(strftime(formatted_dates, "%m"))02 -6.827e-01  7.685e+03   0.000
## as.factor(strftime(formatted_dates, "%m"))03 -1.604e+00  7.427e+03   0.000
## as.factor(strftime(formatted_dates, "%m"))04  1.291e-01  8.144e+03   0.000
## as.factor(strftime(formatted_dates, "%m"))05 -3.020e-02  7.124e+03   0.000
## as.factor(strftime(formatted_dates, "%m"))06 -2.165e+00  6.277e+03   0.000
## as.factor(strftime(formatted_dates, "%m"))07 -2.029e+00  6.182e+03   0.000
## as.factor(strftime(formatted_dates, "%m"))08  1.591e+01  5.904e+03   0.003
## as.factor(strftime(formatted_dates, "%m"))09  1.710e+01  5.904e+03   0.003
## as.factor(strftime(formatted_dates, "%m"))10  1.891e+01  5.904e+03   0.003
## as.factor(strftime(formatted_dates, "%m"))11  1.880e+01  5.904e+03   0.003
## as.factor(strftime(formatted_dates, "%m"))12  7.779e-01  1.043e+04   0.000
##
##               Pr(>|z|)
## (Intercept)         0.9973
## ave_emotion.anger     0.2026
## ave_emotion.anticipation 0.7746
## ave_emotion.disgust    0.5992
## ave_emotion.fear       0.5979
## ave_emotion.joy        0.1149
## ave_emotion.sadness    0.7200

```

```

## ave_emotion.surprise          0.0651 .
## ave_emotion.trust             0.2825
## sentiment_title$ave_sentiment 1.72e-10 ***
## sentiment_text$ave_sentiment  0.0169 *
## review_length                 0.3513
## as.factor(strftime(formatted_dates, "%m"))02 0.9999
## as.factor(strftime(formatted_dates, "%m"))03 0.9998
## as.factor(strftime(formatted_dates, "%m"))04 1.0000
## as.factor(strftime(formatted_dates, "%m"))05 1.0000
## as.factor(strftime(formatted_dates, "%m"))06 0.9997
## as.factor(strftime(formatted_dates, "%m"))07 0.9997
## as.factor(strftime(formatted_dates, "%m"))08 0.9979
## as.factor(strftime(formatted_dates, "%m"))09 0.9977
## as.factor(strftime(formatted_dates, "%m"))10 0.9974
## as.factor(strftime(formatted_dates, "%m"))11 0.9975
## as.factor(strftime(formatted_dates, "%m"))12 0.9999
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 500.40 on 499 degrees of freedom
## Residual deviance: 226.19 on 477 degrees of freedom
## AIC: 272.19
##
## Number of Fisher Scoring iterations: 19
stars_6b <- glm((star_rating_num) > 4 ~ (sentiment_title$ave_sentiment + sentiment_text$ave_sen
, data=reviews_scraped, family = "binomial")
summary(stars_6b)

##
## Call:
## glm(formula = (star_rating_num) > 4 ~ (sentiment_title$ave_sentiment +
## sentiment_text$ave_sentiment + review_length), family = "binomial",
## data = reviews_scraped)
##
## Coefficients:
## Estimate Std. Error z value Pr(>|z|)
## (Intercept) -1.958960 0.238889 -8.200 2.40e-16 ***
## sentiment_title$ave_sentiment 2.801355 0.403621 6.941 3.91e-12 ***
## sentiment_text$ave_sentiment 2.038417 0.528919 3.854 0.000116 ***
## review_length -0.006616 0.002596 -2.548 0.010822 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##

```

```

##      Null deviance: 500.40  on 499  degrees of freedom
## Residual deviance: 353.92  on 496  degrees of freedom
## AIC: 361.92
##
## Number of Fisher Scoring iterations: 6
stars_6_step <- step(stars_6a)

## Start:  AIC=272.19
## (star_rating_num) > 4 ~ (ave_emotion.anger + ave_emotion.anticipation +
##   ave_emotion.disgust + ave_emotion.fear + ave_emotion.joy +
##   ave_emotion.sadness + ave_emotion.surprise + ave_emotion.trust +
##   sentiment_title$ave_sentiment + sentiment_text$ave_sentiment +
##   review_length + as.factor(strftime(formatted_dates, "%m")))
##
##
##              Df Deviance    AIC
## - ave_emotion.anticipation      1    226.28 270.28
## - ave_emotion.sadness            1    226.32 270.32
## - ave_emotion.disgust            1    226.46 270.46
## - ave_emotion.fear               1    226.47 270.47
## - review_length                 1    227.12 271.12
## - ave_emotion.trust             1    227.46 271.46
## <none>                          226.19 272.19
## - ave_emotion.anger             1    228.19 272.19
## - ave_emotion.joy               1    228.86 272.86
## - ave_emotion.surprise          1    230.83 274.83
## - sentiment_text$ave_sentiment   1    232.16 276.16
## - sentiment_title$ave_sentiment  1    285.16 329.16
## - as.factor(strftime(formatted_dates, "%m")) 11    333.02 357.02
##
## Step:  AIC=270.28
## (star_rating_num) > 4 ~ ave_emotion.anger + ave_emotion.disgust +
##   ave_emotion.fear + ave_emotion.joy + ave_emotion.sadness +
##   ave_emotion.surprise + ave_emotion.trust + sentiment_title$ave_sentiment +
##   sentiment_text$ave_sentiment + review_length + as.factor(strftime(formatted_dates,
##   "%m"))
##
##
##              Df Deviance    AIC
## - ave_emotion.sadness            1    226.41 268.41
## - ave_emotion.disgust            1    226.56 268.56
## - ave_emotion.fear               1    226.58 268.58
## - review_length                 1    227.27 269.27
## - ave_emotion.trust             1    227.52 269.52
## <none>                          226.28 270.28
## - ave_emotion.anger             1    228.29 270.29
## - ave_emotion.joy               1    228.95 270.95
## - sentiment_text$ave_sentiment   1    232.22 274.22
## - ave_emotion.surprise          1    233.50 275.50

```

```

## - sentiment_title$ave_sentiment          1    285.20 327.20
## - as.factor(strftime(formatted_dates, "%m")) 11    333.64 355.64
##
## Step: AIC=268.41
## (star_rating_num) > 4 ~ ave_emotion.anger + ave_emotion.disgust +
##     ave_emotion.fear + ave_emotion.joy + ave_emotion.surprise +
##     ave_emotion.trust + sentiment_title$ave_sentiment + sentiment_text$ave_sentiment +
##     review_length + as.factor(strftime(formatted_dates, "%m"))
##
##                                     Df Deviance    AIC
## - ave_emotion.disgust                1    226.58 266.58
## - review_length                      1    227.41 267.41
## - ave_emotion.trust                  1    227.68 267.68
## - ave_emotion.fear                   1    227.85 267.85
## <none>                               226.41 268.41
## - ave_emotion.anger                  1    228.59 268.59
## - ave_emotion.joy                    1    228.98 268.98
## - sentiment_text$ave_sentiment        1    232.60 272.60
## - ave_emotion.surprise                1    233.52 273.52
## - sentiment_title$ave_sentiment       1    285.20 325.20
## - as.factor(strftime(formatted_dates, "%m")) 11    333.79 353.79
##
## Step: AIC=266.58
## (star_rating_num) > 4 ~ ave_emotion.anger + ave_emotion.fear +
##     ave_emotion.joy + ave_emotion.surprise + ave_emotion.trust +
##     sentiment_title$ave_sentiment + sentiment_text$ave_sentiment +
##     review_length + as.factor(strftime(formatted_dates, "%m"))
##
##                                     Df Deviance    AIC
## - review_length                      1    227.54 265.54
## - ave_emotion.trust                  1    227.84 265.84
## - ave_emotion.fear                   1    227.90 265.90
## <none>                               226.58 266.58
## - ave_emotion.anger                  1    228.59 266.59
## - ave_emotion.joy                    1    229.02 267.02
## - sentiment_text$ave_sentiment        1    232.60 270.60
## - ave_emotion.surprise                1    233.52 271.52
## - sentiment_title$ave_sentiment       1    285.40 323.40
## - as.factor(strftime(formatted_dates, "%m")) 11    333.87 351.87
##
## Step: AIC=265.54
## (star_rating_num) > 4 ~ ave_emotion.anger + ave_emotion.fear +
##     ave_emotion.joy + ave_emotion.surprise + ave_emotion.trust +
##     sentiment_title$ave_sentiment + sentiment_text$ave_sentiment +
##     as.factor(strftime(formatted_dates, "%m"))
##
##                                     Df Deviance    AIC
## - ave_emotion.trust                  1    228.82 264.82

```

```

## - ave_emotion.fear          1    229.22 265.22
## - ave_emotion.anger         1    229.48 265.48
## <none>                      227.54 265.54
## - ave_emotion.joy           1    229.98 265.98
## - sentiment_text$ave_sentiment 1    234.07 270.07
## - ave_emotion.surprise       1    234.43 270.43
## - sentiment_title$ave_sentiment 1    287.32 323.32
## - as.factor(strftime(formatted_dates, "%m")) 11    340.65 356.65
##
## Step: AIC=264.82
## (star_rating_num) > 4 ~ ave_emotion.anger + ave_emotion.fear +
##     ave_emotion.joy + ave_emotion.surprise + sentiment_title$ave_sentiment +
##     sentiment_text$ave_sentiment + as.factor(strftime(formatted_dates,
## "%m"))
##
##
##              Df Deviance    AIC
## - ave_emotion.fear          1    230.37 264.37
## <none>                      228.82 264.82
## - ave_emotion.anger         1    230.97 264.97
## - ave_emotion.joy           1    233.31 267.31
## - ave_emotion.surprise       1    234.85 268.85
## - sentiment_text$ave_sentiment 1    235.45 269.45
## - sentiment_title$ave_sentiment 1    288.27 322.27
## - as.factor(strftime(formatted_dates, "%m")) 11    342.80 356.80
##
## Step: AIC=264.37
## (star_rating_num) > 4 ~ ave_emotion.anger + ave_emotion.joy +
##     ave_emotion.surprise + sentiment_title$ave_sentiment + sentiment_text$ave_sentiment +
##     as.factor(strftime(formatted_dates, "%m"))
##
##
##              Df Deviance    AIC
## <none>                      230.37 264.37
## - ave_emotion.anger         1    233.06 265.06
## - ave_emotion.joy           1    235.00 267.00
## - ave_emotion.surprise       1    236.49 268.49
## - sentiment_text$ave_sentiment 1    237.61 269.61
## - sentiment_title$ave_sentiment 1    290.02 322.02
## - as.factor(strftime(formatted_dates, "%m")) 11    348.04 360.04

```

```
summary(stars_6_step)
```

```

##
## Call:
## glm(formula = (star_rating_num) > 4 ~ ave_emotion.anger + ave_emotion.joy +
##     ave_emotion.surprise + sentiment_title$ave_sentiment + sentiment_text$ave_sentiment +
##     as.factor(strftime(formatted_dates, "%m")), family = "binomial",
##     data = reviews_scraped)
##

```

```

## Coefficients:
##
##              Estimate Std. Error z value
## (Intercept)      -2.055e+01  5.941e+03  -0.003
## ave_emotion.anger -1.178e+01  7.907e+00  -1.490
## ave_emotion.joy    1.480e+01  7.791e+00   1.899
## ave_emotion.surprise -1.696e+01  7.834e+00  -2.165
## sentiment_title$ave_sentiment    3.758e+00  5.864e-01   6.409
## sentiment_text$ave_sentiment    1.900e+00  7.268e-01   2.613
## as.factor(strftime(formatted_dates, "%m"))02 -5.850e-01  7.740e+03   0.000
## as.factor(strftime(formatted_dates, "%m"))03 -1.422e+00  7.485e+03   0.000
## as.factor(strftime(formatted_dates, "%m"))04 -1.735e-02  8.132e+03   0.000
## as.factor(strftime(formatted_dates, "%m"))05  4.612e-03  7.147e+03   0.000
## as.factor(strftime(formatted_dates, "%m"))06 -1.746e+00  6.322e+03   0.000
## as.factor(strftime(formatted_dates, "%m"))07 -1.735e+00  6.217e+03   0.000
## as.factor(strftime(formatted_dates, "%m"))08  1.614e+01  5.941e+03   0.003
## as.factor(strftime(formatted_dates, "%m"))09  1.737e+01  5.941e+03   0.003
## as.factor(strftime(formatted_dates, "%m"))10  1.920e+01  5.941e+03   0.003
## as.factor(strftime(formatted_dates, "%m"))11  1.920e+01  5.941e+03   0.003
## as.factor(strftime(formatted_dates, "%m"))12  8.111e-01  1.048e+04   0.000
##
##              Pr(>|z|)
## (Intercept)      0.99724
## ave_emotion.anger    0.13633
## ave_emotion.joy      0.05753 .
## ave_emotion.surprise 0.03042 *
## sentiment_title$ave_sentiment    1.47e-10 ***
## sentiment_text$ave_sentiment    0.00896 **
## as.factor(strftime(formatted_dates, "%m"))02 0.99994
## as.factor(strftime(formatted_dates, "%m"))03 0.99985
## as.factor(strftime(formatted_dates, "%m"))04 1.00000
## as.factor(strftime(formatted_dates, "%m"))05 1.00000
## as.factor(strftime(formatted_dates, "%m"))06 0.99978
## as.factor(strftime(formatted_dates, "%m"))07 0.99978
## as.factor(strftime(formatted_dates, "%m"))08 0.99783
## as.factor(strftime(formatted_dates, "%m"))09 0.99767
## as.factor(strftime(formatted_dates, "%m"))10 0.99742
## as.factor(strftime(formatted_dates, "%m"))11 0.99742
## as.factor(strftime(formatted_dates, "%m"))12 0.99994
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 500.40  on 499  degrees of freedom
## Residual deviance: 230.37  on 483  degrees of freedom
## AIC: 264.37
##
## Number of Fisher Scoring iterations: 19

```



```
AIC(stars_6)
```

```
## [1] 357.0178
```

```
AIC(stars_6a)
```

```
## [1] 272.1933
```

```
AIC(stars_6b)
```

```
## [1] 361.9229
```

```
AIC(stars_6_step)
```

```
## [1] 264.3656
```

Interestingly the full model has the lowest AIC score, even though only the two sentiment scores were significant.

Building a conclusive model for helpfulness:

```
help_a <- glm((N_helpful) > 10 ~ (ave_emotion.anger+ave_emotion.anticipation +
                                ave_emotion.disgust +ave_emotion.fear +
                                ave_emotion.joy +ave_emotion.sadness+
                                ave_emotion.surprise+ave_emotion.trust +
                                sentiment_title$ave_sentiment +
                                sentiment_text$ave_sentiment+review_length),
              data=reviews_scraped, family = "binomial")
summary(help_a)
```

```
##
```

```
## Call:
```

```
## glm(formula = (N_helpful) > 10 ~ (ave_emotion.anger + ave_emotion.anticipation +
##   ave_emotion.disgust + ave_emotion.fear + ave_emotion.joy +
##   ave_emotion.sadness + ave_emotion.surprise + ave_emotion.trust +
##   sentiment_title$ave_sentiment + sentiment_text$ave_sentiment +
##   review_length), family = "binomial", data = reviews_scraped)
##
```

```
## Coefficients:
```

```
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)   -6.269214   0.830195  -7.551 4.30e-14 ***
## ave_emotion.anger    14.943941  14.002050   1.067   0.286
## ave_emotion.anticipation -7.055534  24.519837  -0.288   0.774
## ave_emotion.disgust  -42.019103  36.631222  -1.147   0.251
## ave_emotion.fear     9.403853  32.835695   0.286   0.775
## ave_emotion.joy      1.881173  26.137013   0.072   0.943
## ave_emotion.sadness  42.026397  31.926990   1.316   0.188
## ave_emotion.surprise  4.319313  31.649544   0.136   0.891
## ave_emotion.trust   -15.913280  23.175217  -0.687   0.492
## sentiment_title$ave_sentiment  0.752852  0.747336   1.007   0.314
## sentiment_text$ave_sentiment  1.559535  1.757458   0.887   0.375
```

```
## review_length          0.022442    0.003112    7.212 5.51e-13 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 215.82  on 499  degrees of freedom
## Residual deviance: 108.97  on 488  degrees of freedom
## AIC: 132.97
##
## Number of Fisher Scoring iterations: 8
```

This model to determine if a review will be voted as helpful by more than 10 people shows that only the review_length, not the emotion or sentiment of the review is significant in shaping the outcome.

```
help_a <- glm((N_helpful) > 10 ~ (ave_emotion.anger+ave_emotion.anticipation +
                                ave_emotion.disgust +ave_emotion.fear +
                                ave_emotion.joy +ave_emotion.sadness+
                                ave_emotion.surprise+ave_emotion.trust +
                                sentiment_title$ave_sentiment +
                                sentiment_text$ave_sentiment+review_length+
                                as.factor(strftime(formatted_dates, "%m"))),
              data=reviews_scraped, family = "binomial")
summary(help_a)
```

```
##
## Call:
## glm(formula = (N_helpful) > 10 ~ (ave_emotion.anger + ave_emotion.anticipation +
##   ave_emotion.disgust + ave_emotion.fear + ave_emotion.joy +
##   ave_emotion.sadness + ave_emotion.surprise + ave_emotion.trust +
##   sentiment_title$ave_sentiment + sentiment_text$ave_sentiment +
##   review_length + as.factor(strftime(formatted_dates, "%m"))),
##   family = "binomial", data = reviews_scraped)
##
## Coefficients:
##
##              Estimate Std. Error z value
## (Intercept)    -5.629e+00  1.473e+00  -3.821
## ave_emotion.anger      1.321e+01  2.111e+01   0.626
## ave_emotion.anticipation -1.634e+00  2.650e+01  -0.062
## ave_emotion.disgust    -6.452e+01  4.803e+01  -1.343
## ave_emotion.fear      -9.838e+00  3.862e+01  -0.255
## ave_emotion.joy       -7.844e+00  3.615e+01  -0.217
## ave_emotion.sadness     6.433e+01  3.729e+01   1.725
## ave_emotion.surprise     1.493e+01  3.810e+01   0.392
## ave_emotion.trust     -1.699e+01  2.907e+01  -0.585
## sentiment_title$ave_sentiment  1.628e+00  8.820e-01   1.846
## sentiment_text$ave_sentiment  1.979e+00  2.074e+00   0.954
## review_length         2.622e-02  4.136e-03   6.341
## as.factor(strftime(formatted_dates, "%m"))02 -5.255e-01  1.747e+00  -0.301
```

```

## as.factor(strftime(formatted_dates, "%m"))03 -1.778e+01  1.753e+03  -0.010
## as.factor(strftime(formatted_dates, "%m"))04 -5.793e-01  1.652e+00  -0.351
## as.factor(strftime(formatted_dates, "%m"))05 -1.211e+00  1.486e+00  -0.815
## as.factor(strftime(formatted_dates, "%m"))06 -3.514e-01  1.291e+00  -0.272
## as.factor(strftime(formatted_dates, "%m"))07 -1.544e-01  1.332e+00  -0.116
## as.factor(strftime(formatted_dates, "%m"))08 -1.886e+00  1.384e+00  -1.363
## as.factor(strftime(formatted_dates, "%m"))09 -1.508e+00  1.433e+00  -1.053
## as.factor(strftime(formatted_dates, "%m"))10 -3.593e+00  1.581e+00  -2.272
## as.factor(strftime(formatted_dates, "%m"))11 -2.418e+00  1.606e+00  -1.506
## as.factor(strftime(formatted_dates, "%m"))12  2.215e-01  1.734e+00   0.128
##
##                                     Pr(>|z|)
## (Intercept)                        0.000133 ***
## ave_emotion.anger                  0.531460
## ave_emotion.anticipation           0.950835
## ave_emotion.disgust                0.179178
## ave_emotion.fear                   0.798925
## ave_emotion.joy                    0.828196
## ave_emotion.sadness                0.084458 .
## ave_emotion.surprise               0.695186
## ave_emotion.trust                  0.558839
## sentiment_title$ave_sentiment      0.064888 .
## sentiment_text$ave_sentiment       0.339936
## review_length                     2.29e-10 ***
## as.factor(strftime(formatted_dates, "%m"))02 0.763555
## as.factor(strftime(formatted_dates, "%m"))03 0.991909
## as.factor(strftime(formatted_dates, "%m"))04 0.725796
## as.factor(strftime(formatted_dates, "%m"))05 0.415080
## as.factor(strftime(formatted_dates, "%m"))06 0.785432
## as.factor(strftime(formatted_dates, "%m"))07 0.907724
## as.factor(strftime(formatted_dates, "%m"))08 0.172862
## as.factor(strftime(formatted_dates, "%m"))09 0.292546
## as.factor(strftime(formatted_dates, "%m"))10 0.023071 *
## as.factor(strftime(formatted_dates, "%m"))11 0.132034
## as.factor(strftime(formatted_dates, "%m"))12 0.898316
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##    Null deviance: 215.816  on 499  degrees of freedom
## Residual deviance:  92.566  on 477  degrees of freedom
## AIC: 138.57
##
## Number of Fisher Scoring iterations: 17
help_b <- glm((N_helpful) > 10 ~ (sentiment_title$ave_sentiment +
                                sentiment_text$ave_sentiment+review_length),
              data=reviews_scraped, family = "binomial")

```

```
summary(help_b)
```

```
##
## Call:
## glm(formula = (N_helpful) > 10 ~ (sentiment_title$ave_sentiment +
##     sentiment_text$ave_sentiment + review_length), family = "binomial",
##     data = reviews_scraped)
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)   -5.774236   0.606680  -9.518  < 2e-16 ***
## sentiment_title$ave_sentiment  0.852425   0.711065   1.199   0.231
## sentiment_text$ave_sentiment   0.604435   1.458060   0.415   0.678
## review_length    0.021041   0.002779   7.572 3.68e-14 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##    Null deviance: 215.82  on 499  degrees of freedom
## Residual deviance: 114.42  on 496  degrees of freedom
## AIC: 122.42
##
## Number of Fisher Scoring iterations: 7
```

```
help_c <- glm((N_helpful) > 10 ~ (review_length),
              data=reviews_scraped, family = "binomial")
summary(help_c)
```

```
##
## Call:
## glm(formula = (N_helpful) > 10 ~ (review_length), family = "binomial",
##     data = reviews_scraped)
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)   -5.54935   0.54948 -10.099  < 2e-16 ***
## review_length  0.02013   0.00263   7.654 1.95e-14 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##    Null deviance: 215.82  on 499  degrees of freedom
## Residual deviance: 116.91  on 498  degrees of freedom
## AIC: 120.91
##
## Number of Fisher Scoring iterations: 7
```

```
help_d <- step(help_a)
```

```
## Start:  AIC=138.57
## (N_helpful) > 10 ~ (ave_emotion.anger + ave_emotion.anticipation +
##     ave_emotion.disgust + ave_emotion.fear + ave_emotion.joy +
##     ave_emotion.sadness + ave_emotion.surprise + ave_emotion.trust +
##     sentiment_title$ave_sentiment + sentiment_text$ave_sentiment +
##     review_length + as.factor(strftime(formatted_dates, "%m")))
##
##               Df Deviance    AIC
## - as.factor(strftime(formatted_dates, "%m")) 11  108.969 132.97
## - ave_emotion.anticipation                    1   92.569 136.57
## - ave_emotion.joy                            1   92.613 136.61
## - ave_emotion.fear                          1   92.631 136.63
## - ave_emotion.surprise                      1   92.717 136.72
## - ave_emotion.anger                        1   92.917 136.92
## - ave_emotion.trust                        1   92.929 136.93
## - sentiment_text$ave_sentiment              1   93.460 137.46
## <none>                                     92.566 138.57
## - ave_emotion.disgust                      1   94.627 138.63
## - ave_emotion.sadness                      1   95.365 139.37
## - sentiment_title$ave_sentiment             1   96.247 140.25
## - review_length                           1  189.959 233.96
##
## Step:  AIC=132.97
## (N_helpful) > 10 ~ ave_emotion.anger + ave_emotion.anticipation +
##     ave_emotion.disgust + ave_emotion.fear + ave_emotion.joy +
##     ave_emotion.sadness + ave_emotion.surprise + ave_emotion.trust +
##     sentiment_title$ave_sentiment + sentiment_text$ave_sentiment +
##     review_length
##
##               Df Deviance    AIC
## - ave_emotion.joy                            1   108.97 130.97
## - ave_emotion.surprise                      1   108.99 130.99
## - ave_emotion.fear                          1   109.05 131.05
## - ave_emotion.anticipation                  1   109.06 131.06
## - ave_emotion.trust                        1   109.48 131.48
## - sentiment_text$ave_sentiment              1   109.72 131.72
## - ave_emotion.anger                        1   109.83 131.83
## - sentiment_title$ave_sentiment             1   110.01 132.01
## - ave_emotion.disgust                      1   110.41 132.41
## - ave_emotion.sadness                      1   110.67 132.66
## <none>                                     108.97 132.97
## - review_length                           1   210.55 232.55
##
## Step:  AIC=130.97
## (N_helpful) > 10 ~ ave_emotion.anger + ave_emotion.anticipation +
```

```

##      ave_emotion.disgust + ave_emotion.fear + ave_emotion.sadness +
##      ave_emotion.surprise + ave_emotion.trust + sentiment_title$ave_sentiment +
##      sentiment_text$ave_sentiment + review_length
##
##
##      Df Deviance    AIC
## - ave_emotion.surprise      1    109.00 129.00
## - ave_emotion.fear           1    109.05 129.05
## - ave_emotion.anticipation    1    109.06 129.06
## - ave_emotion.trust          1    109.56 129.56
## - sentiment_text$ave_sentiment 1    109.79 129.79
## - ave_emotion.anger          1    109.87 129.87
## - sentiment_title$ave_sentiment 1    110.03 130.03
## - ave_emotion.disgust        1    110.42 130.42
## - ave_emotion.sadness        1    110.68 130.68
## <none>                      108.97 130.97
## - review_length             1    211.34 231.34
##
## Step:  AIC=129
## (N_helpful) > 10 ~ ave_emotion.anger + ave_emotion.anticipation +
##      ave_emotion.disgust + ave_emotion.fear + ave_emotion.sadness +
##      ave_emotion.trust + sentiment_title$ave_sentiment + sentiment_text$ave_sentiment +
##      review_length
##
##
##      Df Deviance    AIC
## - ave_emotion.anticipation    1    109.06 127.06
## - ave_emotion.fear            1    109.10 127.10
## - ave_emotion.trust           1    109.56 127.56
## - ave_emotion.anger           1    109.88 127.88
## - sentiment_text$ave_sentiment 1    109.89 127.89
## - sentiment_title$ave_sentiment 1    110.04 128.04
## - ave_emotion.disgust         1    110.44 128.44
## - ave_emotion.sadness         1    110.72 128.72
## <none>                       109.00 129.00
## - review_length              1    212.05 230.05
##
## Step:  AIC=127.06
## (N_helpful) > 10 ~ ave_emotion.anger + ave_emotion.disgust +
##      ave_emotion.fear + ave_emotion.sadness + ave_emotion.trust +
##      sentiment_title$ave_sentiment + sentiment_text$ave_sentiment +
##      review_length
##
##
##      Df Deviance    AIC
## - ave_emotion.fear            1    109.14 125.14
## - sentiment_text$ave_sentiment 1    109.92 125.92
## - ave_emotion.anger           1    109.94 125.94
## - ave_emotion.trust           1    110.12 126.12
## - sentiment_title$ave_sentiment 1    110.13 126.13
## - ave_emotion.disgust         1    110.55 126.55

```

```

## - ave_emotion.sadness          1    110.81 126.81
## <none>                          109.06 127.06
## - review_length                1    212.05 228.05
##
## Step: AIC=125.14
## (N_helpful) > 10 ~ ave_emotion.anger + ave_emotion.disgust +
##     ave_emotion.sadness + ave_emotion.trust + sentiment_title$ave_sentiment +
##     sentiment_text$ave_sentiment + review_length
##
##                                Df Deviance    AIC
## - sentiment_text$ave_sentiment  1    109.94 123.94
## - ave_emotion.anger            1    110.02 124.02
## - ave_emotion.trust            1    110.14 124.14
## - sentiment_title$ave_sentiment 1    110.22 124.22
## - ave_emotion.disgust          1    110.64 124.64
## <none>                          109.14 125.14
## - ave_emotion.sadness          1    112.71 126.71
## - review_length                1    212.31 226.31
##
## Step: AIC=123.94
## (N_helpful) > 10 ~ ave_emotion.anger + ave_emotion.disgust +
##     ave_emotion.sadness + ave_emotion.trust + sentiment_title$ave_sentiment +
##     review_length
##
##                                Df Deviance    AIC
## - ave_emotion.trust            1    110.45 122.45
## - ave_emotion.anger            1    110.82 122.82
## <none>                          109.94 123.94
## - ave_emotion.disgust          1    112.00 124.00
## - sentiment_title$ave_sentiment 1    112.03 124.03
## - ave_emotion.sadness          1    113.06 125.06
## - review_length                1    212.41 224.41
##
## Step: AIC=122.45
## (N_helpful) > 10 ~ ave_emotion.anger + ave_emotion.disgust +
##     ave_emotion.sadness + sentiment_title$ave_sentiment + review_length
##
##                                Df Deviance    AIC
## - ave_emotion.anger            1    111.31 121.31
## - sentiment_title$ave_sentiment 1    112.28 122.28
## - ave_emotion.disgust          1    112.41 122.41
## <none>                          110.45 122.45
## - ave_emotion.sadness          1    113.43 123.43
## - review_length                1    213.15 223.15
##
## Step: AIC=121.32
## (N_helpful) > 10 ~ ave_emotion.disgust + ave_emotion.sadness +
##     sentiment_title$ave_sentiment + review_length

```

```
##
##
##           Df Deviance    AIC
## - ave_emotion.disgust      1   112.57 120.57
## - sentiment_title$ave_sentiment 1   113.30 121.30
## <none>                      111.31 121.31
## - ave_emotion.sadness      1   114.52 122.52
## - review_length            1   213.49 221.49
##
## Step:  AIC=120.57
## (N_helpful) > 10 ~ ave_emotion.sadness + sentiment_title$ave_sentiment +
##      review_length
##
##           Df Deviance    AIC
## <none>                      112.57 120.57
## - ave_emotion.sadness      1   114.59 120.59
## - sentiment_title$ave_sentiment 1   115.25 121.25
## - review_length            1   214.48 220.48
```

```
summary(help_d)
```

```
##
## Call:
## glm(formula = (N_helpful) > 10 ~ ave_emotion.sadness + sentiment_title$ave_sentiment +
##      review_length, family = "binomial", data = reviews_scraped)
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)   -6.111571    0.683832  -8.937  < 2e-16 ***
## ave_emotion.sadness    29.272298   18.719725   1.564   0.118
## sentiment_title$ave_sentiment  1.062290    0.664696   1.598   0.110
## review_length     0.021140    0.002787   7.584 3.34e-14 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 215.82  on 499  degrees of freedom
## Residual deviance: 112.57  on 496  degrees of freedom
## AIC: 120.57
##
## Number of Fisher Scoring iterations: 7
```

```
AIC(help_a)
```

```
## [1] 138.5656
```

```
AIC(help_b)
```

```
## [1] 122.4176
```



```
AIC(help_c)
```

```
## [1] 120.909
```

```
AIC(help_d)
```

```
## [1] 120.5742
```